# An Introduction to
# ECONOMETRICS

## Jaydeb Sarkhel
## Santosh Kumar Dutta

# AN INTRODUCTION TO
# ECONOMETRICS

B.A./B.Sc. Economics (Honours)

## Jaydeb Sarkhel

*Retired Professor, Department of Commerce,*
*Burdwan University;*
*Author of 'Microeconomic Theory', 'Macroeconomic Theory', etc.*

## Dr. Santosh Kumar Dutta

*Associate Professor, Department of Economics,*
*Bidhannagar College;*
*Joint Author of 'An Insight into Statistics',*
*'Microeconomics and Statistics' etc.*

**Revised and Enlarged**
**Second Edition**
**2020**

# CBCS (UG) Syllabus on Econometrics for different Universities in West Bengal

## Paper 1.3 | CALCUTTA UNIVERSITY

Economics Core Course X, Core $T_{10}$
—Introductory Econometrics (Sem-IV), FM : 100
(Th : 50 + Tutorial : 30 + Internal assessment : 10 + Attendance : 10)

1. Nature and Scope of Econometrics                    **2 Lecture hours**
1.1 What is Econometrics ?
1.2 Distinction between Economic model and Econometric model
1.3 Concept of Stochstic relation
1.4 Role of random disturbance in econometric model

2. Classical Linear Regression Model (simple linear regression and multiple linear regression) : Part 1                **15 Lecture hours**
2.1 The classical assumptions
2.2 Concepts of population regression function and sample regression function
2.3 Estimation of model by the method of ordinary least squares

3. Classical linear regression model (simple linear regression and multiple linear regression) : Part 2              **15 Lecture hours**
3.1 Properties of Least Square Estimators (BLUE)-Gauss-Markov theorem
3.2 Qualitative (dummy) independent variables (only interpretation of the model)
3.3 Forecasting (only for two variable model) : Expost forecast and Exante forecast

4. Statistical inference in Linear regression model    **20 Lecture hours**
4.1 Sampling distribution of regression estimates : Standard normal, Chi-square, $t$, $F$
4.2 Confidence intervals
4.3 Concepts of Type I and Type II errors
4.4 Testing of hypothesis about $\beta$ and $\sigma^2$ given and with unknown $\sigma^2$ (Standard normal and '$t$' statistics)
4.5 Testing hypothesis involving several parameters : the $F$ test
4.6 Goodness of fit (in terms of $R^2$, adjusted $R^2$ and $F$ statistic)

5. Violations of Classical Assumptions               **10 Lecture hours**
5.1 Multicollinearity-Consequences, Detection and Remedies
5.2 Heteroscedasticity-Consequences, Detection and Remedies
5.3 Autocorrelation-Consequences, Detection and Remedies

6. Specification Analysis                             **10 Lecture hours**
6.1 Omission of a relevant variable
6.2 Inclusion of an irrelevant variable
6.3 Tests of specification errors
6.4 Testing for linearity and normality assumptions.

# CONTENTS

## Chapter 4

## Violations of Classical Assumptions—The Problems of Heteroscedasticity, Autocorrelation and Multicollinearity

**Chapter**

**5**

**Specification Analysis**

# 1

# Definition, Scope and Goals of Econometrics

## 1.1. Definition and Scope of Econometrics

Literally speaking, the word 'econometrics' means "measurement in economics". Econometrics may be considered as the integration of economics, mathematics and statistics for the purpose of providing numerical values for the parameters of economic relationships and verifying economic theories. It is a special type of economic analysis in which the general economic theory formulated in mathematical terms is combined with empirical measurement of economic phenomena. We start from general economic theory, that is, from the relationships of economic variables as suggested by economic theory and express them in mathematical terms. This is called building of an economic model. Next we use statistical methods in order to obtain numerical estimates of the coefficients of the economic relationships. These statistical methods are called *econometric methods.*

Although measurement is an important part of econometrics, the scope of econometrics is much broader, as can be seen from the following quotations :

"Econometrics, the result of a certain outlook on the role of economics, consists of the application of mathematical statistics to economic data to lend empirical support to the models constructed by mathematical economics and to obtain numerical results"[1].

"Econometrics may be defined as the quantitative analysis of actual economic phenomena based on the concurrent development of theory and observation, related by appropriate methods of inference."[2]

"Econometrics may be defined as the social science in which the tools of economic theory, mathematics and statistical inference are applied to the analysis of economic phenomena."[3]

"Econometrics is concerned with the empirical determination of economic laws."[4]

"The art of the econometrician consists in finding the set of assumptions that are both sufficiently specific and sufficiently realistic to allow him to take the best possible advantage of the data available to him."[5]

"Econometricians ... are of positive help in trying to dispel the poor public image of economics (quantitative or otherwise) as a subject in which empty boxes are opened

1. Gerhard Tintner, *Methodology of Mathematical Economics and Econometrics*, The University of Chicago Press, Chicago, 1968, p. 74.
2. P. A. Samuelson, T. C. Koopmans, and J. R. N. Stone, "Report of the Evaluative Committe for Econometrics," *Econometrica*, vol. 22, no. 2, April 1954, pp. 141-146.
3. Arthur S. Goldberger, *Econometric Theory*, John Wiley & Sons, New York, 1964, p. 1.
4. H. Theil, *Principles of Econometrics*, John Wiley & Sons, New York, 1971, P. 1.
5. E. Malinvaud, *Statistical Methods of Econometrics*, Rand Mc Nally, Chicago, 1966, P. 514.

by assuming the existence of can openers to reveal contents which any two economists will interpret in 11 ways."[8]

"The method of econometric research aims, essentially, at a conjunction of economic theory and actual measurements, using the theory and technique of statistical inference as a bridge pier."[7]

All these definitions suggest that econometrics is an amalgam of economic theory, mathematical economics, economic statistics and mathematical statistics.

Economic theory postulates an exact relationship between economic variables but actually an economic relationship always contains a random element. Economic theory ignores it but econometrics does not, because the econometric methods can deal with these random components. For example, in the Keynesian macroeconomic theory we find an exact relationship between consumption expenditure (C) and income (Y). Keynesian consumption function is given by $C = a + bY$ where $a > 0$ is called the

autonomous part of consumption expenditure and $b = \dfrac{dC}{dY}$ is called the marginal

propensity to consume (MPC) [$0 < b < 1$ by assumption]. This is an exact relationship because C is completely determined by Y. So, in this model the effects of other variables like price, wealth, income distribution etc., are ignored. But in econometrics the influence of other factors is considered by introducing a random variable in the model and that random variable is generally denoted by 'u' called error term. So, the consumption function considered in econometrics is $C = a + bY + u$. Now econometrics methods estimate the parameters 'a' and 'b' and while estimating 'a' and 'b' the choice of the econometric method depends on the behaviour of the distribution of the random variable 'u'.

There are three main sources of the error term 'u' in the functional relation. These are :

(i) unpredictable element of randomness in human response,

(ii) effect of a large number of variables that have been omitted from the functional relation,

and (iii) measurement error. (For details see Section 2.2.1)

## 1.2. Relationship between Econometrics and Economic Theory

Economic theory makes statements or hypotheses that are mostly qualitative in nature.

Econometrics presupposes the existence of a body of economic theory. Economic theory should come first because it states the hypothesis about economic behaviour which should be tested with the econometric methods.

For example, we consider the consumption income relationship of the form

$$C = a + bY + u.$$

Economic theory suggests that consumption is a function of income and with the

information of economic theory we know that $MPC = \dfrac{dC}{dY} = b$ lies between 0 and 1 i.e.,

$0 < b < 1$.

6. Adrian C. Darnell and J. Lynne Evans, *The Limits of Econometrics*, Edward Elgar Publishing, Hants, England, 1990, P. 54.

7. T. Haavelmo, "The Probability Approach in Econometrics", Supplement to *Econometrica*, vol. 12, 1944, Preface P. iii.

The proposition suggested by economic theory is to be tested now, applying econometric methods. If we find that the theory is consistent with the empirical results we accept the theory but if we find that it is not consistent with the empirical results, then we have either to reject the theory or to modify the theory. If we like to modify the theory then we should not reject the theory, rather we should incorporate some other variables and parameters to make the theory more meaningful (and close to reality).

For example, the simple consumption income relation, $C = a + bY + u$ can be modified to the form

$C = a + bY + cP + dW + u$ where two new variables $P$ (price level) and $W$ (wealth) have been taken into account in the functional relation. The signs of the parameters $(a, b, c, d > 0)$ and the corresponding response coefficients can also be tested empirically.

### 1.2.1. Difference between Economic Model and Econometric Model

A model is a simplified representation of a real world process. In practice, in any economic model (say consumption function or demand function), we can include all the relevant variables that we think are relevant for our purpose and dump the rest of the variables in a basket called "disturbance". This brings us to the distinction between an economic model and an econometric model.

An economic model is a set of assumptions that approximately describe the behaviour of an economy. An econometric model, on the other hand, consists of the following :

(i) A set of behavioural equations derived from the economic model.

(ii) A statement of whether there are errors of observation in the observed variables.

(iii) A specification of the probability distribution of the "disturbances" (and errors of measurement).

For example, we may consider a simple demand model of economics. Then econometric model will usually consist of :

(a) The behavioural equation : $q = \alpha + \beta p + u$ where $q$ = quantity demanded, $p$ = price, $\alpha$ and $\beta$ are two parameters and $u$ = random disturbance term.

(b) A specification of probability distribution of $u$, where values of $u$ are independently and normally distributed with mean $E(u) = 0$ and variance $(u) = \sigma_u^2$. With these specifications we can test empirically the law of demand or the hypothesis that $\beta < 0$.

We may also use the estimated demand function for prediction and policy purposes.

## 1.3. Econometrics and Mathematical Economics

Mathematical economics states economic theory in terms of mathematical symbols. There is no essential difference between mathematical economics and economic theory. Both state the same relationships, but while economic theory uses verbal exposition, mathematical economics employs mathematical symbolism. Both express the various economic relationships in an exact form. Neither economic theory nor mathematical economics allows for random elements which might affect the relationship and make it *stochastic*. Furthermore, they do not provide numerical values for the coefficients of the relationships. Relations in economic theory or in mathematical economics are of *non-stochastic* form. It is in this regard that econometrics differs from mathematical economics.

Although econometrics presupposes the expression of economic relationships in mathematical form, like mathematical economics it does not assume that economic relationships are exact. On the contrary, econometrics assumes that relationships are not exact. Econometric methods are designed to take into account random disturbances which create deviations from the exact behavioural patterns suggested by economic theory and mathematical economics. Furthermore, econometric methods provide numerical values of the coefficients of economic phenomena. Thus, by combining mathematical formulations of theory with empirical data, econometrics enables us to pass from the abstract theoretical scheme to numerical results in concrete cases.

## 1.4 Econometrics and Statistics

Econometrics differs both from mathematical statistics and economic statistics. An economic statistician gathers empirical data, records them, tabulates them or charts them and then attempts to describe the pattern in their development and perhaps detect some relationship between various economic magnitudes. Thus economic statistics is mainly a descriptive aspect of economic theory. It does not provide explanations of the development of the various variables and does not provide measurement of the parameters of economic relationships.

Economic statistics differs from mathematical or inferential statistics. Mathematical statistics is based upon the theory of probability and deals with the methods of measurement which are developed on the basis of controlled or carefully planned experiments. These statistical methods cannot be applied to economic relationships because such experiments cannot be designed except in a very few cases, e.g. agricultural experiments or industrial experimentation for economic phenomena.

Econometrics uses statistical methods after adopting them to the problems of economic life. These adopted statistical methods are called econometric methods. In particular, econometric methods are so adjusted that they become appropriate for the measurement of economic relationships which are stochastic, that is, they include random elements. The adjustments consists primarily in specifying the stochastic random elements that are supposed to operate in the real world and enter into the determination of the observed data so that the latter can be interpreted as a random sample to which the methods of statistics can be applied.

## 1.5. Goals of Econometrics

Econometrics helps us to achieve the following three main goals:

(i) **Analysis**: This means testing of economic theories. There are alternative theories to explain the functioning of the economic system. Econometrics examines the explanatory power of the system.

(ii) **Policy making**: The numerical estimates of the coefficients of the economic relationships help the policy-maker to define the appropriate policies. For example, the numerical estimate of price elasticities of demand for a product will help the policy maker to know how much additional revenue is expected to be obtained if sales tax is imposed on that commodity. Alternatively, numerical estimates of price elasticities of exports and imports will help us to know how far the devaluation as a policy will be effective in solving the balance of payments deficit problem.

(iii) **Forecasting**: The numerical estimates of the coefficients are used in order to forecast the future value of the economic variables. Without forecasting the planner cannot adopt appropriate policies. Of course, these goals are not mutually exclusive

Successful econometric applications should really include some combination of all those skills.

## 1.6 Division of Econometrics

Econometrics may be divided into two branches — theoretical econometrics and applied econometrics

*Theoretical econometrics* includes the development of appropriate methods for the measurement of economic relationships. Econometric techniques are based on statistical techniques which have been adopted to the particular characteristics of economic relationships.

Econometric methods may be classified into two groups (i) Single equation techniques which are methods that are applied to one relationship at a time and (ii) Simultaneous equation techniques, which are methods applied to all the relationships of a model simultaneously.

*Applied econometrics* includes the applications of econometric methods to specific branches of economic theory. It examines the problems encountered and the findings of applied research in the fields of demand, supply, production, investment, consumption and other sectors of economic theory. Applied econometrics involves the application of the tools of theoretical econometrics for the analysis of economic phenomena and forecasting economic behaviour.

## 1.7. Methodology/Stages of Econometric Research

Applied econometric research is concerned with the measurement of the parameters of economic relationships and with the prediction of the values of economic variables.

The relationships of economic theory which can be measured with one or another econometric techniques are causal, that is they are relationships in which some variables are postulated as causes of the variation of other variables. In this sense definitional equations do not require any measurement but examine the equation $Y = C + I$ the mathematical expression of the definition of national income in a closed economy with no government activity of economic theory. It does not explain the determination of the level of income or the causes of its variations.

There are four stages in any econometric research

### Stage A — Specification of the model .

It means expressing the relationships between the variables in mathematical form. This stage is also called formulation of the maintained hypothesis involves the determination of

(i) dependent and the explanatory variables to be included in the model

(ii) the theoretical expectations about the sign, size of the parameters of the function

(iii) the mathematical form of the model

For example, consider a production function of the following type $Y = f(K, L)$ where $K$ and $L$ are the two factors of production.

[$K$ = Capital, $L$ = Labour and $Y$ is the level of output] This function can also be written in the Cobb-Douglas form i.e $Y = K^\alpha L^\beta$ or $\log Y = \alpha \log K + \beta \log L$. This is the mathematical form of log linear function. Here some theoretical restrictions must be imposed $0 < \alpha, \beta < 1,$

$\alpha + \beta > 1$ if there are increasing returns to scale

$\alpha + \beta < 1$ if there are decreasing returns to scale

$\alpha + \beta = 1$ if there are constant returns to scale

$\alpha$ = Elasticity of output with respect to capital

$\beta$ = Elasticity of output with respect to labour

## Stage B   Estimation of the model

[several lines illegible]

... the variables of the function
... among the explanatory variables
... $\beta_1 P + \beta_2 Y + \beta_3 W + u$ where $t$
... endogenous variables and explanatory variables are ...
... Of course we have to find out whether there is ...
... among the parameters variable in the equation law or ...
... the problem of multicollinearity.)

... the appropriate econometric technique for the estimation of the
function and critical examination of the assumptions of the chosen technique and of
their possible implications for the estimates of the coefficients.

## Stage C   Evaluation of estimates

After the estimation of the model the econometrician must proceed with the
evaluation of the results of the calculations, that is with the determination of the
reliability of these results. The evaluation consists of deciding whether the estimates
of the parameters are theoretically meaningful and statistically significant.

For this purpose we may use various criteria which may be classified into three
groups.

i.  **Economic criteria**   These are determined by the principles of economic theory
and refer to the sign and the size of the parameters of economic relationships. For
example the Keynesian liquidity preference function may be expressed in the
mathematical form

$$M = \beta_0 + \beta_1 Y + \beta_2 r + u$$

where $M$ = demand for money (dependent variable), $Y$ = income, $r$ = rate of interest,
$u$ = error term, $\beta_0$, $\beta_1$, $\beta_2$ are the parameters whose values and signs are to be
determined on the basis of observed data. On the basis of the existing theory the signs
of the parameters would be $\beta_0 > 0$, $\beta_1 > 0$, $\beta_2 < 0$.

ii.  **Statistical criteria (First order tests)**   These are determined by statistical theory
and aim at the evaluation of the statistical reliability of the estimates of the parameters
of the model. The most widely used statistical criteria are the correlation coefficient and
the standard error of the estimates.

iii.  **Econometric criteria (Second order tests)**   These are set by the theory of
econometrics and aim at the investigation of whether the assumptions of the
econometric method employed are satisfied or not in any particular case. The
econometric criteria serve as second order tests (as tests of the statistical tests). In
other words they determine the reliability of the statistical criteria, and in particular

the standard errors of the parameter estimate. This helps to establish whether the estimates ... in a plausible manner ... using ... t-statistics ...

**Stage 9  Evaluation of the forecasting power of the estimated model**

... accepting the model ...

... test the forecasting power of the model.

## 1.8 Desirable Properties of an Econometric Model

... of an econometric model ... the ... following desirable properties

i) **Theoretical plausibility**. The model ... compatible with the postulates ... economic theory. It must describe the economic phenomenon to which it relates.

ii) **Explanatory ability**. The model should be able to explain the observations of the actual world. It must be consistent with the observed behaviour of the economic variables whose relationship it determines.

iii) **Accuracy of the estimates of the parameters**. The estimates of the coefficients should be accurate in the sense that they should approximate as best as possible the true parameters of the structural model.

iv) **Forecasting ability**. The model should produce satisfactory predictions of future values of the dependent variables.

v) **Simplicity**. The model should represent the economic relationships with maximum simplicity.

## 1.9 Nature and Sources of Data for Economic Analysis

The success of any econometric analysis depends on the availability of the appropriate data. Three types of data are generally available for empirical analysis: *time series data*, *cross section data* and *pooled data/panel data*.

### Time Series Data

A time series is a set of observations on the values that a variable takes at different times. Such data may be collected at regular time intervals, such as *daily* (e.g. stock prices, weather reports), *weekly* (e.g. money supply figures), *monthly* (e.g. unemployment rate, Consumer Price Index (CPI)), *quarterly* (e.g. GDP), *annually* (e.g. government budget), *quinquennially* (that is every 5 years (e.g. the census of manufactures) or *decennially* (that is every 10 years (e.g. the census of population)).

### Cross-Section Data

Cross-section data are data on one or more variables collected at the same point of time, such as the census of population conducted by the Government of India every 10 years, the Survey of household consumer expenditure in India conducted by National Sample Survey Organization (NSSO), the opinion polls by the Times of India, NDTV, CNN-IBN and many other organizations. An individual researcher or a group may also collect cross-section data directly from the field of enquiry.

Conventionally, the letter $Y$ denotes the dependent variable and $X$'s $(X_1, X_2, \ldots, X_p)$ denote the explanatory independent variables, $X_k$ being the $k$th explanatory variable. The subscript $i$ or $t$ denote $i$th or $t$th observation or value $X_k$ or $Y_k$ will denote the $i$th (or $t$th) observation on variable $Y_k$. Here $N$ (or $T$) will denote the total number of observations or values in the population and $n$ (or $t$) will denote the total number of observations in a sample. Normally the subscript $i$ will be used for cross-section data (i.e. data collected at one point of time) and the subscript $t$ will be used for time series data (i.e. data collected on different periods of time). For instance, consider the Keynesian consumption function of the form $C = a + bY$ where $C$ = consumption expenditure, $Y$ = income and $a$ and $b$ are two constants, $a$ = autonomous part of consumption expenditure, $b$ = marginal propensity to consume. According to the existing theory $a > 0, 0 < b < 1$. If we like to test this relation with the help of time series data then we will write the regression equation in the form $C = a + bY + u_t$ (where $t = 1, 2, \ldots, t$ (say)) where $u$ is the random disturbance term. On the other hand, we can write the regression equation in the form $C = a + bY + u_i$, $i = 1, 2, \ldots, N$ (say), when we verify the relation with the help of cross-section data.

## Pooled Data

Pooled or combined data are elements of both time series and cross section data. Generally speaking, pooled data is a combination of data (i.e. sales, advertisements, earnings etc.) of say 20 firms over a given period of time say a year or two. These combined data of 20 firms in 2 years making 40 observations make is a pooled data, that is, pooling 20 firms data in 2 years together. So, it is a combination of cross section data and time series data.

## Panel, Longitudinal, or Micropanel Data

This is a special type of pooled data in which the same cross-sectional unit (say a family or a firm) is surveyed over time.

For example, the Nigerian population commission surveys each house every 10 years to determine the changes that may have occurred within these years. By surveying or interviewing the same households or firms to find out their population or financial conditions periodically (10 years interval), panel data can help to provide useful information on the changes that may have occurred in these households. It is more detailed than just the pooled data in a short period of time.

Thus there is a basic difference between pooled data and panel data. It should be noted that pooled time-series, cross-section data are data with relatively few cross-sections (say few firms under study), where variables are held in cross-section specific individual series (i.e. sales, advertisement, earnings, etc.), while panel data correspond to data with large number of cross-sections, with variables held in single series in stacked form.

## The Sources of Data

The data used in empirical analysis may be collected by a government agency (e.g. the Central Statistical Organization), an international agency (e.g., the International Monetary Fund (IMF) or the World Bank), a private organization (e.g. the Centre for Monitoring Indian Economy) or an individual. There exist a lot of agencies collecting data for one purpose or another. Nowadays the Internet has revolutionized data

gathering. Most of the data can be downloaded from different websites either free of cost or with minimum cost.

### The Accuracy of Data

Although plenty of data are available for economic research, the quality of data is often not that good.

There are several reasons for that:

(i) Most of the social science data are non experimental in nature. Therefore there is the possibility of observational errors.

(ii) Even in experimentally collected data, errors of measurement arise from approximations and rounding offs.

(iii) In questionnaire type of surveys, the problem of non-response may lead to bias in results.

(iv) The sampling methods used in obtaining data may vary so widely that it is often difficult to compare the results obtained from the various samples.

(v) Economic data are generally available at a highly aggregate level. Such highly aggregated data may not be helpful for individualistic study.

Because of all of these and many other problems, the researchers should always keep in mind that the results of research are only as good as the quality of the data. Therefore, if in given situations researchers find that the results of the research are "unsatisfactory" the cause may not be that they used the wrong model, but due to the poor quality of data.

## 1.10. A Note on the Measurement Scales of Variables

The variables that we generally use can be measured in four types of scales : *ratio scale, interval scale, ordinal scale* and *nominal scale*. We can briefly describe them as follows

**Ratio Scale** : For a variable $X$, taking two values say $X_1$ and $X_2$, the ratio $X_1/X_2$ and the distance $(X_2 - X_1)$ are meaningful quantities. Also, there is a natural ordering (ascending or descending) of the values along the scale (say $X_2 \geq X_1$ or $X_1 \leq X_2$). Most economic variables belong to this category. Personal income, measured in rupees is a ratio variable, someone earning ₹ 50,000 is making twice as much as another person earning ₹ 25000.

**Interval scale** : The interval scale satisfies the last two properties stated in ratio scale but not the first.

For example, the distance between two time periods, say (2018, 2001) is meaningful

but not the ratio of two time periods $\left(\dfrac{2018}{2001}\right)$

**Ordinal Scale** : A variable belongs to this category only if it satisfies the third property of the ratio scale (i.e. natural ordering). Examples are grading systems (A, B, C, grades) or income class (upper, middle, lower). For these variables the ordering exists. But the distances between the categories cannot be quantified.

**Nominal Scale** : Variables in this category have more of the features of the ratio scale variables. Variables such as gender (male, female) and marital status (married, unmarried, divorced, separated) simply denote categories.

## EXERCISE

1. What is Econometrics and what are its components? Describe the functions of each component with examples in support of your answer.

2. Name and describe three relationships studied in Economic theory which can be examined as subject matter of Econometrics. What are the parameters of these relationships?

3. How would you define Econometrics? How does it differ from Mathematical Economics and Statistics? Describe the main steps involved in any econometric research by taking an example from economic theory.

4. Considering the following relations, how would you explain that economic theory postulates exact relationships between economic variables? How can these be transformed into econometric relations?

   Demand function $D = a + b P$, $b < 1$ when $D$ = quantity demanded, $P$ = price and

   i. income

   ii. Supply function $S = a + b P$ where $S$ = quantity supplied, $P$ = price

   iii. Consumption function $C = a + b Y_d$, where $C$ = consumption expenditure and $Y_d$ = disposable income

   iv. Cost function $= a + b t$ where $t$ = total cost and $k$ = total output

   v. Production function $t = a L^a K^b$ where $y$ = level of output, $L$ = labour input, $K$ = capital input, $A$ = constant technical parameter, $a$, $b$ are the two elasticity coefficients

   a. What is the economic meaning of the coefficients involved in all the above equations?

   b. What would you expect about the sign and size of the coefficients in each of the above relationships?

5. Enumerate the relation between econometrics and economic theory.

6. What is Econometrics? What are the different goals of econometrics?

7. Distinguish between theoretical econometrics and applied econometrics.

8. Explain briefly the different stages of any econometric research.

9. What is an econometric model? Illustrate any one of such models.

10. What are the desirable properties of an econometric model?

11. What are the different types of data available for empirical analysis?

12. Distinguish between time series data and cross-section data.

13. What are pooled data? What are panel data? Distinguish between pooled data and panel data.

14. What are the different sources of data used in empirical analysis?

15. What do you mean by accuracy of data? What are the different reasons for distortion of accuracy of data collected and published by different organizations?

16. Give a brief outline on measurement scales of variables.

# 2

# The Simple Linear Regression Model

## 2.1. Introduction

Most of economics is concerned with relations among variables. These relations when phrased in mathematical terms can predict the effect of one variable on another. For example, assuming that income, prices of other commodities and all other determinants of demand are constants. We can express the quantity demanded $q$ of any commodity as a function of the price $(p)$ of that commodity only. This may be put in the form $q = f(p)$. Similarly we are familiar with other functions with different assumptions such as consumption function $C = f(t)$, supply function $S = f(p)$, cost function $C = f(q)$, production function $Q = f(x_1, x_2)$ where $x_1$ and $x_2$ are amounts of different inputs, etc.

These functional relationships define the dependence of the dependent variable upon the independent variable(s) in the specific form. The functional relation may be linear, quadratic, logarithmic, exponential or hyperbolic.

A relation between two variables X and Y expressed as $Y = f(X)$ is said to be deterministic or non-stochastic (non-random) if for each value of the independent variable (X) there is one and only one corresponding value of the dependent variable (Y). On the other hand, a relation between X and Y is said to be stochastic if for a particular value of X there is a whole probability distribution of values of Y. In such a case for any given value of X, the dependent variable Y assumes some specific value only with some probability.

For example, a linear demand function (in deterministic form) can be written as $q = f(p) = \alpha - \beta p$ ($\alpha > 0$, $\beta < 0$) and in particular $q = 100 - 5p$. When $p = 10$, $q = 50$, when $p = 15$, $q = 25$ etc.

But such an exact and deterministic relation between $p$ and $q$ is never true in the real world.

The deterministic behaviour of the above relationship breaks down when the *ceteris paribus* (other things remaining the same) condition is relaxed.

We therefore rewrite the demand equation as $q = \alpha - \beta p + u$ or in particular $q = 100 - 5p + u$ where $u$ is commonly known as *random disturbance* since it disturbs an otherwise deterministic relation.

### 2.1.1. Concepts of Population Regression Function and Sample Regression Function

Sampling denotes the selection of a part of the aggregate statistical material with a view to obtaining information about the whole. This aggregate or totality of statistical information on a particular character of all the members covered by an investigation is called population or universe. When the population size is very large it may not be

11

possible to take a complete enumeration of the population. Then we select a small part of the population and by examining it we can infer about the nature of the whole population. The basic idea of sampling is to make inference about the population by examining a small part of it.

In practice we may be interested to find out the relation between two or more variables simultaneously. In the case of simple (linear) regression model we assume only one explanatory variable but in the case of multiple regression model we assume more than one explanatory variables. The first case is known as the bivariate analysis while the second case is known as the multivariate analysis. In this chapter we will concentrate on bivariate analysis (study the relation between two variables X and Y where Y is dependent variable, X independent explanatory variable).

We know that regression analysis is largely concerned with estimating and/or predicting the population parameter say mean value of the dependent variable (Y) on the basis of the known or fixed values of the explanatory variables. To understand the fact we consider a total population of 60 families in a hypothetical community and their monthly income (X) and monthly consumption expenditure (Y), both in rupees. These 60 families are divided into 10 income groups and the monthly expenditures of each family in the various groups are shown in the following table (Table 2.1).

**Table 2.1** Joint distribution of monthly income (X in ₹) and monthly consumption expenditure (Y in ₹) of 60 families in a hypothetical community

| X→ Y↓ | 8000 | 10000 | 12000 | 14000 | 16000 | 18000 | 20000 | 22000 | 24000 | 26000 |
|---|---|---|---|---|---|---|---|---|---|---|
| | 5500 | 6500 | 7900 | 8000 | 10200 | 11000 | 12000 | 13500 | 13700 | 15000 |
| | 6000 | 7000 | 8400 | 9300 | 10700 | 11500 | 13000 | 14700 | 14500 | 15200 |
| | 6500 | 7400 | 9000 | 9500 | 11000 | 12000 | 14000 | 4000 | 15500 | 17500 |
| | 7000 | 8000 | 9400 | 10400 | 11600 | 13000 | 4400 | 5200 | 16500 | 7800 |
| | 7500 | 8500 | 9800 | 10900 | 11900 | 13500 | 4500 | 15700 | 17500 | 8000 |
| | | 8800 | | 11300 | 12500 | 14000 | | 16000 | 18900 | 8500 |
| | | | | 11500 | | | | 16200 | | 9000 |
| Total | 32500 | 46200 | 44500 | 70700 | 67800 | 75000 | 68500 | 104500 | 96500 | 121000 |
| Conditional means of Y E Y/X) | 6500 | 7700 | 8900 | 10100 | 11300 | 12500 | 13700 | 14900 | 16100 | 17300 |

Here we have 10 fixed values of X and the corresponding Y values against each of the X values and hence we have 10 subpopulations of Y. From Table 2.1 we see that there is considerable variation in monthly consumption expenditure in each income group but the general picture is that despite the variability of monthly consumption expenditure within each income bracket, on an average monthly consumption expenditure increases as income increases. To understand it clearly we have given the mean, or average monthly consumption expenditure corresponding to each of the 10 levels of income. Thus, corresponding to the monthly income level of ₹8000, the mean consumption expenditure is ₹6500 and so on. In total we have 10 mean values for 10 sub-populations of Y and these mean values are called conditional expected values as they depend upon the given values of the (conditioning) variable X.

Symbolically we denote them as $E(Y \mid X)$ which simply means the expected value of $Y$ given the value of $X$. It should be noted that these expected values are called conditional expected values. In order to calculate the conditional expected values $E(Y \mid X)$ we have to construct conditional probability distribution of $Y$, $P(Y \mid X)$, shown in Table 2.2.

**Table 2.2   Conditional Probabilities $P(Y \mid X_i)$ for the data of Table 2.1**

| $X \rightarrow$ $P(Y \mid X_i)$ | 800 | 1000 | 1200 | 1400 | 1600 | 1800 | 2000 | 2200 | 2400 | 2600 |
|---|---|---|---|---|---|---|---|---|---|---|
| Conditional probabilities $P(Y \mid X_i)$ | $\frac{1}{5}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{7}$ |
| | $\frac{1}{5}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{7}$ |
| | $\frac{1}{5}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{7}$ |
| | $\frac{1}{5}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{7}$ |
| | $\frac{1}{5}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | $\frac{1}{5}$ | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{7}$ |
| | | $\frac{1}{6}$ | | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{6}$ | | $\frac{1}{7}$ | $\frac{1}{6}$ | $\frac{1}{7}$ |
| | | | | $\frac{1}{7}$ | | | | $\frac{1}{7}$ | | $\frac{1}{7}$ |
| Conditional means of $Y$ | 6500 | 7700 | 8900 | 10100 | 11300 | 12500 | 13700 | 14900 | 16100 | 17300 |

For the income group of ₹800 the expected monthly expenditure is obtained by

$$₹5500 \times \frac{1}{5} + ₹6000 \times \frac{1}{5} + ₹6500 \times \frac{1}{5} + ₹7000 \times \frac{1}{5} + ₹7500 \times \frac{1}{5} = ₹6500$$

The expected monthly expenditures for other income groups are also obtained in this way.

It is important to distinguish these conditional expected values from the *unconditional expected value* of monthly consumption expenditure $E(Y)$. If we add the monthly consumption expenditures for all the 60 families in the population and divide this number by 60, we get the value ₹12120 (₹727200÷60) which is the unconditional mean or expected value of $Y$, $E(Y)$.

Thus the expected monthly consumption expenditure of a family would be ₹12120 (the unconditional mean). But if we like to know the expected value of monthly consumption expenditure of a family whose monthly income is say ₹1200, then we get a value of ₹8900 (the conditional mean).

Graphically, if we join these conditional mean values, we obtain the **population regression line (PRL)** or **population regression curve** or simply it is the *regression of Y on X*.

The population regression curve is simply the locus of the conditional means of the dependent variable for the fixed values of the explanatory variable(s). More specifically,

It is the curve which may be meant as the obsorvations of Y corresponding to the given values of the regressor X. This is shown in Figure 2.2
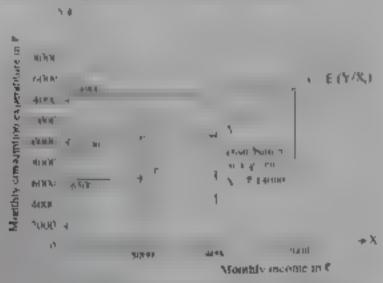


Fig. 2.2. Population regression line, data of Table 2.1

This figure shows that for each X (i.e. income level) there is a population of Y values (monthly consumption expenditure) that are spread around the conditional mean of those Y values. For simplicity we have assumed that these Y values are distributed symmetrically around their respective (conditional) mean values and the regression line (or curve) passes through these (conditional) mean values.

## 2.1.2. Population Regression Function (PRF)

From the above explanation and Figure 2.1 it is clear that each conditional mean $E(Y|X)$ is a function of $X_i$, where $X$ is a given value of X. Symbolically $E(Y|X) = f(X_i)$ where $f(X_i)$ denotes some function of the explanatory variable X, a linear function of $X_i$. This function is also known as the **conditional expectation function (CEF)** or **Population regression function (PRF)** or **population regression (PR)**. It states that the expected value of the distribution of Y given $X_i$ is functionally related to $X_i$. In simple terms, it tells how the mean or average response of Y varies with X. However, the functional form of the PRF is an empirical question. For example, to verify the consumption-income relation we generally assume a linear relation. We may assume that the PRF $E(Y|X_i)$ is a linear function of $X_i$, say of the type

$E(Y|X_i) = \alpha + \beta X$ where $\alpha$ and $\beta$ are unknown but fixed parameters known as the regression coefficients. $\alpha$ and $\beta$ are also known as intercept and slope coefficients respectively. In regression analysis our interest is in estimating the PRF and the unknown values of $\alpha$ and $\beta$ on the basis of observations on Y and X.

From our earlier example stated in Table 2.1 we see that, given the income level of $X$ (say), an individual family's consumption expenditure is clustered around the average consumption of all families at $X_i$, i.e. around its conditional expectation. Therefore, we can express the deviation of an individual $Y_i$ around its expected value as follows

$u_i = Y - E(Y|X)$ or $Y = E(Y|X)$ or $Y = \alpha + \beta X + u$

where the deviation $u_i$ is an unobservable random variable taking positive or negative values. Technically, $u_i$ is known as the **Stochastic disturbance** term or **Stochastic error term**.

### 2.1.3 The Sample Regression Function (SRF)

In practice, what we have to do is to estimate the population regression function (PRF) on the basis of the sample information. For example, we assume that the population is not known to us. But it is real. Now we have a sample of 9 values of Y and X as shown in Table 2.3. In that table, we have the corresponding to the given X's, each Y given in Table 2.1. Later, assuming we have a sample which is the same Table 2.1 as the population values.

| Table 2.3 | |
|---|---|
| A random sample from the population of Table 2.1 | |
| Y | X |
| 7000 | 80000 |
| 6500 | 60000 |
| 9000 | 71000 |
| 9500 | 81000 |
| 1000 | 76000 |
| 1500 | 81000 |
| 12000 | 200000 |
| 14000 | 220000 |
| 15500 | 240000 |
| 13000 | 260000 |

| Table 2.4 | |
|---|---|
| A random sample from the population of Table 2.1 | |
| Y | X |
| 7500 | 80000 |
| 6000 | 80000 |
| 8000 | 100000 |
| 9000 | 120000 |
| 1500 | 160000 |
| 9000 | 180000 |
| 8400 | 200000 |
| 15500 | 220000 |
| 9500 | 240000 |
| 8500 | 260000 |

Now from the sample of Table 2.3 we can predict or forecast the average monthly consumption expenditure Y in the population as a whole corresponding to the chosen X's.

However, we may not be able to estimate the PRF accurately because of sampling fluctuations. To see this, we have drawn another random sample from the same population (Table 2.1) and shown in Table 2.4.

Now plotting the data of Tables 2.3 and 2.4 we obtain the scatter diagram shown in Figure 2.2.

In the scatter diagram two sample regression lines are drawn so as to fit the scatters reasonably well. $SRF_1$ is based on the first sample and $SRF_2$ is based on the second sample. The regression lines in Figure 2.2 are known as the sample regression lines. Supposedly they represent the population regression line but due to sampling fluctuations they are at best an approximation of the true PRF. In general, we may get N different SRF's for N different samples and these SRF's are not likely to be same.

Like the PRF (derived from population regression line) we can develop the concept of the **sample regression function (SRF)** to represent the sample regression line.

Since from the population regression function we know that the function is of the form, $Y = E(Y|X_i) = \alpha + \beta X$, the sample counterpart of this equation may be written as

$$\hat{Y} = \hat{\alpha} + \hat{\beta} X, \text{ where}$$

$\hat{Y}$ = estimator of $E(Y|X_i)$, $\hat{\alpha}$ = estimator of $\alpha$ and $\hat{\beta}$ = estimator of $\beta$

It should be noted that an **estimator** also known as a (sample) *statistic* is a rule or formula or method that tells how to estimate the population parameter from the information provided by the sample at hand. A particular numerical value obtained by the estimator is known as an **estimate**. It should be noted that an estimator is random but an estimate is non random.



**Fig. 2.2.** Regression lines based on two different samples

Like the population regression function $(Y_i = E(Y/X_i) = u_i = \alpha + \beta X_i + u_i)$ we can express the sample regression equation $Y_i = \alpha + \beta X_i$ in its stochastic form as follows

$Y_i = \alpha + \beta X_i + w_i$ where $w_i$ denotes the (sample) residual term. Conceptually $w_i$ is analogous to $u_i$ and can be regarded as an estimate of $u_i$. It is introduced in the *SRF* for the same reasons as $u_i$ was introduced in the *PRF*.

Now our primary objective in regression analysis is to estimate the *PRF* $Y_i = \alpha + \beta X_i + u_i$ on the basis of the *SRF* $Y_i = \alpha + \beta X_i + w_i$. It should be noted that the estimate of the *PRF* based on the *SRF* is at best an approximate one, as sampling fluctuations exist). We have to develop procedures that tell us how to construct the *SRF* in that the *PRF* as faithfully as possible.

## 2.2. The Simple Linear Regression Model

Relationships suggested by economic theory are usually specified as exact or deterministic relationships between variables while on the other hand much stress is placed on the need for testing these economic theories. This implies a belief in the existence of stochastic function. The knowledge of econometrics tries to test these theoretical propositions in terms of stochastic variables. The simplest form of stochastic relation between two variables $X$ and $Y$ is called a simple linear regression model and is given by $Y_i = \alpha + \beta X_i + u_i$ for $i = 1, 2, \ldots n$ where $Y$ = dependent variable $X$ explanatory variable (independent variable), $u$ = stochastic disturbance term, $\alpha$ and $\beta$ are two regression parameters whose values are to be determined on the basis of the

given data on $X$ and $Y$. The subscript $i$ refers to the $i$th observation, $n$ sample size or number of data points.

The stochastic nature of the regression model implies that for every value of $X$ there is a whole probability distribution of values of $Y$. In other words the value of $Y$ can never be predicted exactly. This uncertainty concerning $Y$ is due to the presence of the stochastic term $u$ which imparts randomness to $Y$.

### 2.2.1 Role of Random disturbance term in Econometric Model

We may ask why should we add an error term or random disturbance term $u$ to an econometric model ?

The disturbance term $u$ is a surrogate for all those variables that are omitted from the model but that collectively affect $Y$. We may give the following reasons for the insertion of the disturbance term in an econometric model.

i **Omission of variables from the function** Suppose in the model $Y = \alpha + \beta X + u$ the variable $Y$ denotes the consumption expenditure and $X$ denotes disposable income. But in reality $X$ is not the only variable influencing $Y$. The family's size, tastes the family's spending habits and so on affect the variable $Y$. The error $u$ is a which for the effects of all these variables, some of which may not even be quantifiable and some of which may not even be identifiable. Therefore $u$ may be used as a substitute for all the excluded or omitted variables from the model.

(ii) **Unpredictable element of randomness in human responses** For instance if $Y$ consumption expenditure of a household and $X$ disposable income of the household there is an unpredictable element of randomness in each household's consumption. The household does not behave like a machine. In one month the people in the household are on a spending spree in other month they are tightfisted.

iii **Imperfect specification of the mathematical form of the model** We may have linearised a possibly nonlinear relationship between $X$ and $Y$ or we may have left out of the model some equations.

It is because the economic phenomena are much more complex than a single equation may reveal. For example price determines and is determined by the quantity supplied (or quantity demanded) in the market. Under such circumstances if we attempt to study the phenomena with a single equation model, we are bound to commit an error. Thus the disturbance term represents such an error which may be due to imperfect specification of the form of the model, that is, of the number of equations.

iv **Core variables versus peripheral variables** In consumption income relation for instance we may observe that besides income $X_1$ the number of children per family $X_2$, sex $X_3$, region $X_4$, education $X_5$, and geographical region $X_6$ etc. also can affect consumption expenditure. But it is quite possible that the joint influence of all or some of these variables may be so small that as a practical matter it does not pay to introduce them into the model explicitly. However their combined effect can be treated as a random variable $u_i$.

(v) **Principle of Parsimony** Generally we would like to keep our regression model as simple as possible. If we can explain the behaviour of $Y$ substantially with two or three explanatory variables and if our theory is not strong enough to suggest what other variables might be included, why introduce more variables ? In such cases $u$

... Due to aggregation ...

... Due to errors in measurement (Errors in variables) ...

... errors in the observations.

For all these reasons, the stochastic disturbances is assumed an extremely ... role in regression analysis.

## 2.5 Classical Linear Regression Model and Its Assumptions

Let us consider an observed relation between two variables $X$ and $Y$ which is given by $Y_i = \alpha + \beta X_i + u_i$ for $i = 1, 2, \ldots, n$ where $Y$ is the dependent variable, $X$ is the independent variable, $u$ is the disturbance term, the subscript $i$ denotes the $i$th $i$th observations etc. and $\alpha$ and $\beta$ are the two parameters whose values are to be estimated on the basis of the observed data on $X$ and $Y$.

Now the model is called a Classical Linear Regression Model (CLRM) if the model satisfies the following properties/assumptions.

**Assumption 1** $u$ is a random variable which follows normal distribution.

**Assumption 2** $E(u_i) = 0$ for each $i = 1, 2, \ldots, n$. This means that the probability distribution of the disturbance term is such that its mean is zero.

Now $E(u_i) = 0$ implies $E(Y_i) = \alpha + \beta X_i$. This can be shown as follows. Since $Y_i = \alpha + \beta X_i + u_i$. Now $E(Y_i) = E(\alpha + \beta X_i + u_i) = E(\alpha) + \beta E(X_i) + E(u_i) = \alpha + \beta X_i$ as $E(u_i) = 0$ and $E(X_i) = X_i$.

But $\alpha + \beta X_i$ is the true value of $Y$. This means that expectation of observed value of the dependent variable is its true value. In other words the probability distribution of $Y$ is centred around the true relationship.

**Assumption 3** Variance of each $u_i$ is a constant and is independent of $i = 1, 2, \ldots, n$ and is denoted by $\sigma_u^2$ or simply $\sigma^2$

i.e. $Var(u_i) = \sigma_u^2$ or $\sigma^2$

or $E[u_i - E(u_i)]^2 = E(u_i)^2 = \sigma_u^2$ where $E(u_i) = 0$

Assumptions 2 and 3 imply that the random variables $u_1, u_2, \ldots, u_n$ are identically distributed with the same mean (zero) and same variance ($\sigma_u^2$)

i.e. $u_i \sim IID(0, \sigma_u^2)$ for each $i = 1, 2, 3, \ldots n$

**Assumption 4** The different error terms are independent, distributed $\sim (u_i, u_j)$ ...

Now $\text{Cov}(u_i, u_j) = 0$ for $i \neq j$

... $\sigma_u^2$ ... where ...

**Assumption 5** The independent variable $X$ is non-stochastic ... an input variable and is uncorrelated with explanatory variables

$X_i ... $ for $i = 1 ... n$

The regression equation $Y = \alpha + \beta X$ ... along with the ... assumptions represents the Classical Linear Regression Model. The last assumption is important ... roles to play in the sampling distributions of parameters ...

The role of first three assumptions on the probability distribution of dependent variable $Y$ can now be rationalised.

(i) In the equation $Y = \alpha + \beta X + u$, $Y$ is a linear function of $u$, since $u$ is normally distributed it follows that $Y$ is also normally distributed.

(ii) $Y_i = \alpha + \beta X_i + u_i$

$E(Y_i) = E(\alpha + \beta X_i + u_i)$, $\qquad E(u_i) = \alpha + \beta X_i$

$= \alpha + \beta X_i$

This means that the mean of $Y_i$ is $\alpha + \beta X_i$

(iii) $\text{Var}(Y_i) = E[Y_i - E(Y_i)]^2 = E[Y_i - E(Y_i)]^2 = E[\alpha + \beta X_i + u_i - \alpha - \beta X_i]^2$

$$= E(u_i)^2 = \sigma_u^2 \left[ \because E(u_i)^2 = \sigma_u^2 \right]$$

Therefore, we say that variance of $Y_i$ is $\sigma_u^2$

Thus with the first three assumptions of $u$, we can directly say that $Y$ is normally distributed with mean $(\alpha + \beta X_i)$ and variance $\sigma_u^2$

Symbolically, $Y_i \sim N(\alpha + \beta X_i, \sigma_u^2)$ when $u_i \sim N(0, \sigma_u)$. This is illustrated in Fig. 2.3.

Let $Y = \alpha + \beta X$ represent the population regression line. This regression line is unknown as we do not know the exact values of $\alpha$ and $\beta$. We have to estimate the values $\alpha$ and $\beta$ on the basis of sample data.



Fig. 2.3.

20

[text too faded to read reliably]



**Fig. 2.4.**

in Fig 2.4, AB is the true regression line, CD is the estimated regression and represents one of the observations in the sample data, e differs from u because the true values of the parameters are different from their estimated values. In fact, we can think the residual e as the estimate of the disturbance $u_i$.

## 2.4 Methods of Estimating Regression Parameters

There are different methods for estimating the regression parameters.

Now we shall discuss three methods for estimating the regression parameters $\alpha$ and $\beta$. These are

The method of moments

(ii) The method of ordinary least squares (OLS)

(iii) The method of maximum likelihood (MLE)

## 2.5. The Method of Moments

The assumptions we have made (Assumption 4) about the error term $u$ imply that $E(u) = 0$ and $\text{Cov}(X, u) = 0$.

In the method of moments, we replace these conditions by their sample counterparts

Let $\hat{\alpha}$ and $\hat{\beta}$ be the estimators of $\alpha$ and $\beta$ respectively.

Since $Y_i = \alpha + \beta X_i + u_i$ is the regression equation.

The sample counterpart of $u_i$ is the estimated error term which is also termed the residual defined as $\hat{u}_i = Y_i - \hat{\alpha} - \hat{\beta} X_i$.

The two equations to determine $\alpha$ and $\beta$ are obtained by replacing population assumptions by their sample counterparts.

| Population Assumption | Sample Counterpart |
|---|---|
| $E(u_i) = 0$ | $\sum_{i=1}^{n} \hat{u}_i = 0$ (or) $\sum \hat{u}_i = 0$ |
| $\text{cov}(X_i, u_i) = 0$ | $\sum_{i=1}^{n} X_i \hat{u}_i = 0$ (or) $\sum X_i \hat{u}_i = 0$ |

Thus we get the two equations

$$\sum_{i=1}^{n} \hat{u}_i = 0 \text{ or } \sum_{i=1}^{n} (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0$$

and

$$\sum_{i=1}^{n} X_i \hat{u}_i = 0 \text{ or } \sum_{i=1}^{n} X_i (Y_i - \hat{\alpha} - \hat{\beta} X_i) = 0$$

These equations can be written as

$$\sum_{i=1}^{n} Y_i = n\hat{\alpha} + \hat{\beta} \sum X_i \qquad \text{......(1)}$$

$$\sum_{i=1}^{n} X_i Y_i = \hat{\alpha} \sum X_i + \hat{\beta} \sum X_i^2 \qquad \text{...(2)}$$

These two equations are called "normal equations". Solving these two equations we can get $\alpha$ and $\beta$.

**Example 2.1.** Consider the data on advertising expenditures $(X)$ and sales revenue $(Y)$ for an athletic sports wear store for 5 months.

The observations are as follows

| Month | Sales Revenue $(Y)$ (in 000 ₹) | Advertising Expenditure $(X)$ (in 00 ₹) |
|---|---|---|
| 1 | 2 | 1 |
| 2 | 4 | 2 |
| 3 | 2 | 3 |
| 4 | 6 | 4 |
| 5 | 8 | 5 |

**Solution** Let $Y_i = \alpha + \beta X_i + u_i$ be the regression equation. The two normal equations for estimating the regression coefficients are

$$\sum_{i=1}^{n} Y = n\hat{\alpha} + \hat{\beta} \sum_{i=1}^{n} X_i, \qquad (1)$$

$$\sum_{i=1}^{n} X_i Y = \hat{\alpha} \sum_{i=1}^{n} X_i + \hat{\beta} \sum_{i=1}^{n} X_i^2 \qquad (2)$$

Calculations for $\alpha$ and $\beta$ (estimators of $\alpha$ & $\beta$)

| Month | x | y | xy | $x^2$ | ŷ |
|---|---|---|---|---|---|
|  | 1 | 3 | 4 | 3 | .5 |
| 2 | 2 | 4 | 4 | 8 | 1.6 |
|  | 3 | 5 | 4 | 6 | 6 |
|  | 3 | 5 |  | 24 | 9.2 |
| 4 | 4 | 6 | 36 | 8 |  |
| 5 | 5 | 8 | 8 |  |  |

Total   $\sum x = 5$   $\sum y = 3$   $\sum x = 3$   $\sum x^2 = 81$   $\sum \hat{y} =$

where $n = 5$. Now $\bar{x}$ in the two normal equations

we get   $5\alpha + 15\beta = $   (1)

$15\alpha + 55\beta = 8$   (2)

Solving (1) and (2) by Cramer's rule we get

$$\alpha = \frac{\begin{vmatrix} 8 & 15 \\ 5 & 4 \end{vmatrix}}{\begin{vmatrix} 5 & 15 \\ 15 & 55 \end{vmatrix}} = \frac{2\times5 \quad 155 \quad 50}{3+4 \quad 5 \quad 50} \qquad \text{and} \quad \beta = \frac{\begin{vmatrix} 5 & 15 \\ 15 & 8 \end{vmatrix}}{\begin{vmatrix} 5 & 15 \\ 15 & 55 \end{vmatrix}} = \frac{15 \quad 5}{5 \quad 5} = \frac{1}{1}$$

Thus the estimated regression equation is

$\hat{y} = \beta_1 + \beta_2 x = $

The intercept $\hat{\alpha}$ gives the value of $\hat{y}$ when $x = 0$, this says that advertising expenditures are zero, sales revenue will be $\hat{y}$ million. The slope coefficient is... says that if an advertisement expenditure (X) is changed by one unit, $\hat{y}$ increases (sales) distance by ... unit and $\hat{y}$ ... on an average. We have seen that the estimated value of the residuals given by

$\hat{u}_i = y_i - \hat{y}_i$, shown in the last column of the above table.

## 2.6. The Method of Ordinary Least Squares (OLS)

Let $y = \alpha + \beta x + u$ be a two-variable linear regression model where $y$ is the dependent variable and $x$ is the independent variable and $u$ is the disturbance term. If the disturbance term $u$ satisfies the following properties, then this model will be a classical linear regression model (CLRM)

$E(u_i) = 0$ for each $i = 1, 2, 3, \ldots n$

ii)  $Var(u_i) = \sigma_u^2$ for each $i$

iii)  $E(u_i u_j) = 0$ for all $i \neq j$

iv)  $E(u_i u_j) = \sigma_u^2$ for all $i = j$

v)  Independent variable $X$ is non-stochastic

The two parameters $\alpha$ and $\beta$ of the regression equation can be obtained by the method of ordinary least squares (OLS). Let $\alpha$ and $\beta$ be the estimated values of

and $\beta$. The estimated relation becomes $Y = \alpha + \beta X$, and $e$ = ... is the error term, which shows the difference between the observed and estimated value.

The method of least squares consists in finding out those values of $\alpha$ and $\beta$ for which $\sum_{i=1}^{n} e_i^2$ is minimum. This means that we have to minimize $\sum_{i=1}^{n} e^2 = \sum_{i=1}^{n}$ ...

$\sum_{i=1}^{n} (Y_i - \alpha + \beta X_i)$ through the choice of $\alpha$ and $\beta$. The necessary conditions of minimization require

$$\frac{\delta \sum_{i=1}^{n} e_i^2}{\delta \alpha} = 2 \sum_{i=1}^{n} (Y_i - \alpha + \beta X_i) = 0 \qquad (1)$$

and

$$\frac{\delta \sum_{i=1}^{n} e^2}{\delta \beta} = 2 \sum_{i=1}^{n} X_i (Y_i - \alpha + \beta X_i) = 0 \qquad (2)$$

Simplifying equations (1) and (2) we get two normal equations

$$\sum_{i=1}^{n} Y = n\alpha + \beta \sum_{i=1}^{n} X \qquad (3)$$

$$\sum_{i=1}^{n} XY = \alpha \sum_{i=1}^{n} X + \beta \sum_{i=1}^{n} X_i \qquad (4)$$

Now solving equations (3) and (4) by Cramer's rule we have

$$\beta = \frac{\begin{vmatrix} n & \sum_{i=1}^{n} Y \\ \sum_{i=1}^{n} X & \sum_{i=1}^{n} X_i Y_i \end{vmatrix}}{\begin{vmatrix} n & \sum_{i=1}^{n} X \\ \sum_{i=1}^{n} X & \sum_{i=1}^{n} X^2 \end{vmatrix}} = \frac{n \sum_{i=1}^{n} XY - \sum_{i=1}^{n} X \sum_{i=1}^{n} Y}{n \sum_{i=1}^{n} X^2 - \sum_{i=1}^{n} X}$$

or $\beta = \dfrac{\sum_{i=1}^{n} x_i y}{\sum_{i=1}^{n} x_i^2}$ assuming $X$, $X - \bar{X}$ and $Y = Y - \bar{Y}$

Again from equation (3) we get

$$\sum_{i=1}^{n} Y_i = m\alpha + \beta \sum_{i=1}^{n} X_i$$

$$\sum_{i=1}^{n} \frac{Y_i}{n} = \alpha + \beta \sum_{i=1}^{n} \frac{X_i}{n} \quad \text{or} \quad \bar{Y} = \alpha + \beta \bar{X} \qquad \alpha = \bar{Y} - \beta \bar{X}$$

### 2.6.1 Reverse Regression

By applying O.L.S method we have estimated the linear regression equation $Y = \alpha + \beta X + u$ where $u$ satisfies all the properties of CLRM. The estimated regression equation becomes $\hat{Y} = \hat{\alpha} + \hat{\beta} X$ where $\hat{\alpha}$ and $\hat{\beta}$ are the OLS estimators of $\alpha$ and $\beta$.

Here $\hat{\beta} = \dfrac{\sum_{i=1}^{n} X_i Y_i - \bar{X} \bar{Y}}{\sum \frac{X_i^2}{n} - \bar{X}^2} = \dfrac{cov(X,Y)}{var(X)} = \dfrac{\sigma_{XY}}{\sigma_X^2} = r_{XY} \dfrac{\sigma_Y}{\sigma_X}$

if we put $x = X - \bar{X}$ and $y = Y - \bar{Y}$ then we have

$$\hat{\beta} = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2} \quad \text{and} \quad \hat{\alpha} = \bar{Y} - \hat{\beta} \bar{X}$$

Here $\hat{\beta}$ is the estimated regression coefficient of $Y$ on $X$. In this case the regression equation so defined is called the direct regression equation (of $Y$ on $X$) where $Y$ is the dependent variable and $X$ is the independent variable. Sometimes we have to consider the regression equation of $X$ on $Y$ as well. This is called **reverse regression.**

The reverse regression is used in many cases. For instance reverse regression has been advocated in the analysis of sex (or race) discrimination in salaries.

Suppose $Y$ = salary and $X$ = qualification and we are interested in determining if there is sex discrimination in salaries. We can ask

1. Whether men and women with the same qualifications (value of $X$) are getting the same salaries (value of $Y$). This question is answered by the direct regression, i.e. regression of $Y$ on $X$. Alternatively, we can ask
2. Whether men and women with same salaries (value of $Y$) have the same qualifications (value of $X$).

This question is answered by the reverse regression, i.e. regression equation of $X$ on $Y$.

For the reverse regression, the regression equation can be written as $X = \alpha' + \beta' Y + v$, where $v$ are the errors satisfying all the properties of CLRM. Here $X$ is the dependent variable and $Y$ is the independent variable.

The estimated relation becomes $Y = \alpha + \beta X$ and $r = Y - \hat{Y}$ which is the residual term showing the difference between the observed and estimated values. The method of least squares consists in finding out those values of $\alpha$ and $\beta$ for which

$\sum r^2$ is minimum. This means that we have to minimize

$$\sum r^2 = \sum (Y - \hat{Y})^2 = \sum (Y - \alpha - \beta X)^2 \qquad (1)$$

through the choice of $\alpha$ and $\beta$. The necessary conditions of minimization require

$$\frac{\partial \sum r^2}{\partial \alpha} = -2\sum (Y - \alpha - \beta Y) = 0 \qquad \qquad \qquad 1$$

and

$$\frac{\partial \sum r^2}{\partial \beta} = -2\sum Y (X - \alpha - \beta Y_i) = 0 \qquad \qquad 2$$

Simplifying equations (1) and (2) we get two normal equations.

$$\sum_{i=1}^{n} X_i = n\alpha + \beta \sum_{i=1}^{n} Y \qquad \qquad \qquad 3$$

$$\sum_{i=1}^{n} X_i Y_i = \alpha \sum_{i=1}^{n} Y_i + \beta \sum_{i=1}^{n} Y_i^2 \qquad \qquad 4$$

Now solving equations (3) and (4) by Cramer's rule we have,

$$\beta = \frac{\begin{vmatrix} n & \sum_{i=1}^{n} X_i \\ \sum_{i=1}^{n} Y_i & \sum_{i=1}^{n} X_i Y_i \\ n & \sum_{i=1}^{n} Y_i \\ \sum_{i=1}^{n} Y_i & \sum_{i=1}^{n} Y^2 \end{vmatrix}}{} = \frac{n\sum_{i=1}^{n} X_i Y_i - \sum_{i=1}^{n} X_i \sum_{i=1}^{n} Y_i}{n\sum_{i=1}^{n} Y_i^2 - \left(\sum_{i=1}^{n} Y_i\right)^2}$$

$$= \frac{Cov(X, Y)}{Var(Y)} = \frac{r_{XY}\sigma_X \sigma_Y}{\sigma_Y^2} = r_{XY}\frac{\sigma_X}{\sigma_Y}$$

$$= \text{Regression coefficient of } X \text{ on } Y$$

If we ... ... then the ... express ... follow

$$\beta = \frac{\sum \cdots}{\sum \cdots}$$

Again ... regression ... just

$$\sum \cdots = m \cdots + \sum \cdots$$

or $\quad \sum \frac{1}{n} \cdots = \beta \sum \frac{1}{n} \cdots$ ... $\quad \beta \bar{Y}$

It should be noted that $\beta$ is the regression coefficient of $Y$ on $X$ and ... regression coefficient of $X$ on $Y$.

Since $\beta = r_{xy} \frac{\sigma_y}{\sigma_x}$ and $\beta' = r_{xy} \frac{\sigma_x}{\sigma_y}$

$$\beta \beta' = r_{xy} \frac{\sigma_y}{\sigma_x} \cdot r_{xy} \frac{\sigma_x}{\sigma_y} = r_{xy}^2$$

The two regression lines $Y$ on $X$ and $X$ on $Y$ will be different
... The two regression lines will coincide if $r_{xy} = \pm 1$ and they will
perpendicular to each other if $r_{xy} = 0$.

**Example 2.11.** We now consider a numerical example where we want ... both the
regression direct regression (i.e. $Y$ on $X$) and reverse regression ($X$ on $Y$) on
the basis of the following data

| X | 10 | 1 | 10 | 5 | 8 | 8 | 6 | 7 | 9 | b |
|---|----|---|----|---|---|---|---|---|---|---|
| Y | 1 | 10 | 2 | 6 | 10 | 7 | 9 | 10 | | 11 |

where $X$ = labour-hours of work, $Y$ = output

**Solution.** We know that the fitted direct regression equation $Y$ on $X$ is given ...

$$Y = \alpha + \beta X \quad \text{where } \beta = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$$

and $\alpha = \bar{Y} - \beta \bar{X}$ where $x_i = X_i - \bar{X}$, $Y$ and $y_i = Y_i - \bar{Y}$. Conversely, the equation of
the fitted reverse regression equation ($X$ on $Y$) is given by,

$$\bar{X} = \alpha - \beta Y \quad \text{where } \beta = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} y_i^2}, \quad \text{where } \alpha = \bar{X} - \beta \bar{Y} \text{ and } x_i = X_i - \bar{X}$$

Calculations for direct and reverse regression lines

| | | | | | | $\bar{y}$ | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | 11 | | | | 4 | 4 | 4 | |
| | | 7 | 10 | -1 | 0.4 | 0.4 | 1 | | |
| | | 10 | 12 | 2 | 2.4 | 4.8 | 4 | | 5.76 |
| | | | | | | | | | |
| | | 8 | 10 | 0 | 0.4 | 0 | 8 | | 0.16 |
| | | 11 | 9 | 0 | 2.6 | 0 | 8 | | 6.76 |
| | 1 | | 9 | | | | | | |
| | 8 | 7 | 10 | 1 | 4 | 4 | | | |
| | 9 | 9 | 11 | 1 | 1.4 | 1.4 | 1 | | 1.96 |
| | 10 | 11 | 10 | | 4 | 8 | 4 | | |
| Total | | | 11 | | | | | | |
| | | 80 | 96 | | | | | | 0.4 |

$$\bar{x} = \frac{80}{n} = \frac{80}{1} \quad \bar{y} = \frac{13}{n} = \frac{96}{10} \quad 9.6$$

Now $\quad \beta = \dfrac{\sum xy}{\sum x^2} = \dfrac{21}{28} \quad 0.75$

$$\hat{\alpha} = \bar{y} - \beta\bar{x} = 9.6 \quad 0.75 \quad \bar{x} \quad \cdot 6$$

Estimated regression equation of $Y$ on $X$ (direct regression) is given by
$Y = \alpha + \beta x$ or, $\hat{y} = 3.6 + 0.75 X$

Again, $\quad \gamma = \dfrac{\sum xy}{\sum y^2} \quad \dfrac{21}{40.4} \quad 0.690$

and $\alpha$ $\quad \bar{x} - \gamma \bar{y} = \bar{x} \quad 0.690 \times 9.6 = 1.376$

Estimated reverse regression equation ($X$ on $Y$) is given by $x = \alpha + \beta y$
or, $\quad X = 1.376 + 0.690 Y$

It should be noted that $\beta = 0.75$ is the estimated regression coefficient of $Y$ on $X$
and $\beta = 0.690$ is the estimated regression coefficient of $X$ on $Y$.

Since $\beta \times \beta = r_{XY}^2$

$r_{XY}^2 = 0.75 \times 0.690 \quad 0.5175 = 0.52$ and $r_{XY} \quad \sqrt{0.52} \quad 0.72$

## 2.5.2. Scaling and Units of Measurement

In the regression analysis the units in which the regressand or the dependent variable $Y$ and the regressor(s) are measured make difference in the regression results.

Suppose we like to regress Indian gross domestic savings (GDS) and gross domestic product (GDP), in rupees crore as well in rupees lakh measured in 1999, 2000 prices. We

a researcher assume that in the regression of ADS on GDP one researcher uses data in rupees crore and another researcher uses data in rupees lakh. Now the natural question will the regression result be the same in both cases? The data units in which the regressand and regressors are measured make any difference in the regression result? What is the sensible course to follow in choosing units of measurement for regressand and regressor? To answer these questions, let us proceed as follows.

Let
$$Y = \alpha + \beta X + u_i \qquad (1)$$

where $Y$ = ADS and $X$ = GDP. Let us define

$$Y^* = W \cdot Y \qquad (2)$$

and
$$X^* = W \cdot X \qquad (3)$$

where $W$ and $W'$ are constants, called the scale factors. $W$ may be equal to $W'$ or may be different. If $X$ and $Y$ are measured in say rupees crore and we want to express them in rupees lakh, we will have $Y^* = 100 \, Y$ and $X^* = 100 \, X$ here $W = W' = W$

Now consider the regression using $Y^*$ and $X^*$ variables

$$Y^* = \alpha^* + \beta^* X^* + u_i^* \qquad (4)$$

where $Y^* = W \, Y$, $X^* = W' \, X$ and $u_i^* = W \, u_i$, or $W' = W_2$

Now comparing equations (1) and (4) we can find out the relationships between the following pairs.

1. $\alpha$ and $\alpha^*$

2. $\beta$ and $\beta^*$

3. Var $(\alpha)$ and Var $(\alpha^*)$

4. Var $(\beta)$ and Var $(\beta^*)$

5. $\sigma_u^2$ and $\sigma_u^{2*}$

6. $r_{XY}^2$ and $r_{X^*Y^*}^2$

From the least squares theory we know that [applying OLS method on equation]

$$\alpha = \bar{Y} - \beta \bar{X} \qquad (5)$$

$$\beta = \frac{\Sigma x_i y_i}{\Sigma x_i^2} \qquad (6) \text{ where } x_i = X - \bar{X} \text{ and } y = Y - \bar{Y}$$

$$\text{Var } (\alpha) = \frac{\Sigma X_i^2}{n \Sigma x_i^2} \, \sigma_u^2 \qquad (7)$$

$$\text{Var } (\beta) = \frac{\sigma_u^2}{\Sigma x_i^2} \qquad (8)$$

and $$\sigma_u^2 = \frac{\Sigma u_i^2}{n-2} \text{ or } \frac{\Sigma u_i^2}{n-2} \qquad (9)$$

Similarly, applying OLS method to equation 4, we obtain

$$\hat{\alpha}^* = \bar{Y}^* - \hat{\beta}^* \bar{X}^* \qquad \text{(10)}$$

$$\hat{\beta}^* = \frac{\ }{\ } \qquad \text{where } y^*, \ x^*, \ \text{and} \ y^* \ldots$$

$$\text{Var } \hat{\alpha}^* = \frac{\sum x^{*2}}{n \sum} \cdot \eta_u^* \qquad \text{(12)}$$

$$\text{Var } \hat{\beta}^* = \frac{\eta_u^{*2}}{\sum} \qquad \text{(13)}$$

$$\eta_u^{*2} = \frac{\sum \hat{u}^{*2}}{n-2} \quad \text{or} \quad \frac{\sum \hat{e}^{*2}}{n-2} \qquad \text{(14)}$$

Thus we see that from model (1) $\hat{\alpha}$ and $\hat{\beta}$ are the OLS estimators of $\alpha$ and $\beta$ and from model (2) $\alpha^*$ and $\beta^*$ are the OLS estimators of $\alpha^*$ and $\beta^*$. From the above results it is easy to establish relationship between two sets of parameters.

Since $Y^* = W_1 Y$ (or $y_i^* = w_1 y_i$), $X_i^* = W_2 X_i$ (or $x_i^* = w_2 x_i$), $\hat{u}^* = w_1 u_i$, $\bar{Y}^* = W_1$ and $\bar{X}^* = W_2 \bar{X}$ we can easily verify that

$$\hat{\beta}_i^* = \left( \frac{W_1}{W_2} \right) \hat{\beta} \qquad \text{(15)}$$

$$\hat{\alpha}^* = W_1 \hat{\alpha} \qquad \text{(16)}$$

$$\eta_u^{*2} = W_1^2 \eta_u^2 \qquad \text{(17)}$$

$$\text{Var } (\hat{\alpha}^*) = W_1^2 \text{ Var } (\hat{\alpha}) \qquad \text{(18)}$$

$$\text{Var } (\beta^*) = \frac{W_1}{W_2}^2 \text{ Var } (\beta) \qquad \text{(19)}$$

$$r_{X^*Y^*}^2 = r_{XY}^2 \qquad \text{(20)}$$

From the above results it is clear that from the regression results based on one scale of measurement, we can derive the results based on another scale of measurement once the scaling factors are known. From the results given in (15) to (20) we can also derive some special cases. For instance, if the scaling factors are identical (i.e. $W_1 = W_2$), the slope coefficient and its standard error remain unaffected in going from the $(Y, X)$ to the $(Y^*, X^*)$ scale. However, the intercept and its standard error are both multiplied by $W$ (when $W = W_1$). But if the $X$ scale is not changed (i.e. $W_2 = 1$) and the $Y$ scale is changed by the factor $W$, the slope as well as the intercept coefficients and their respective standard errors are all multiplied by the same $W$ factor. Finally, if the $Y$ scale

remains unchanged ... but the ... scale is changed by the ... If ... the slope coefficient and its standard error are multiplied by the factor ... but the ... intercept coefficient and the standard error remain unaffected.

It should, however, be noted that the transformation in ... variables from ... to ... the $(Y, X)$ scale does not affect the properties of OLS estimators.

To illustrate the above theoretical results we consider an example showing the relationship between GDS and GDP in India during the period 1955-56 to 1984-85.

The estimated regression equation of GDS on GDP both GDS and GDP in rupees crore is given by

$$\hat{GDS}_t = 6742.57 + 0.16 GDP_t \qquad (2.1)$$
SE              (772.00) (0.02)   $r^2 = 0.8891$

Similarly, the estimated regression equation of GDS on GDP both GDS and GDP in rupees lakh is given by

$$\hat{GDS}_t = 674214.57 + 0.16 GDP_t \qquad (2.2)$$
SE              (77210.0 74) (0.02)   $r^2 = 0.8891$

Here we see that the intercept and its standard error is 100 times the corresponding values in the regression (2.1) we should note that $\beta = 100$ (going from crore to lakhs of rupees, i.e., crore = 100 lakhs), but the slope coefficient as well as its standard error is unchanged, in accordance with the theory.

Now suppose we measure GDS in rupees crore and GDP in rupees lakh, the estimated regression equation becomes

$$\hat{GDS}_t = 6742.57 + 0.0016 GDP_t \qquad (2.3)$$
SE              (772.00) (0.0002)   $r^2 = 0.8891$

As expected, the slope coefficient as well as the standard error is $\frac{1}{100}$ th value of equation (2.1) since only the GDP scale is changed.

If we express GDS in rupees lakh and GDP in rupees crore, the estimated regression equation becomes

$$\hat{GDS}_t = 674276.57 + 16.33 GDP_t \qquad (2.4)$$
SE              (77200.74) (1.78)   $r^2 = 0.8891$

Here we see that both the intercept and the slope coefficients as well as their respective standard errors are 100 times their values in equation (2.1) in accordance with our theoretical results.

It should be noted that the $r^2$ value remains the same in all the cases as it is invariant to changes in the unit of measurement and scales.

## 2.7 Estimation of a Function whose Intercept is Zero

In some cases economic theory postulates relationships which have a zero intercept that is, they pass through the origin of the $XY$ plane. [For example, long run consumption function of the form $C = bY$ where $b = APC = MPC$, $C =$ consumption expenditure $=$ income].

In this event we should estimate the function $y = \alpha + \beta x + u$, imposing the

restriction $\alpha = 0$. The formula for the estimation of $\beta$ then becomes $\beta = \dfrac{\sum\limits_{i=1}^{n} x_i y_i}{\sum\limits_{i=1}^{n} x_i^2}$ which

involves the actual values of the variables, and not their deviations, as in the case of unrestricted value of $\alpha$.

**Proof** We want to fit the line $y = \alpha + \beta x + u_i$, subject to the restriction $\alpha = 0$. To estimate $\beta$, the problem is put in a form of restricted minimization problem and then Lagrange method is applied.

Now we have to minimize

$$\sum_{i}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - \alpha - \beta x_i)^2$$

subject to $\alpha = 0$

The Lagrange composite function then becomes $L = \sum_{i}^{n} (y_i - \alpha - \beta x_i)^2 + \lambda \alpha$ where

$\lambda$ is the Lagrange multiplier. Now we have to minimize $L$ with respect to $\alpha$, $\beta$ and $\lambda$. First order conditions of minimization require

$$\frac{\partial L}{\partial \alpha} = -2\sum_{i}^{n} (y_i - \alpha - \beta x_i) + \lambda = 0 \qquad (1)$$

$$\frac{\partial L}{\partial \beta} = -2\sum_{i}^{n} (y_i - \alpha - \beta x_i)x_i = 0 \qquad (2)$$

$$\frac{\partial L}{\partial \lambda} = \alpha = 0 \qquad (3)$$

Now substituting (3) in (2) and rearranging we get

$$-2\sum_{i=1}^{n} x_i (y_i - \beta x_i) = 0$$

or $\sum_{i=1}^{n} x_i y_i - \beta \sum_{i}^{n} x_i^2 = 0$ $\qquad \beta = \dfrac{\sum_{i}^{n} x_i y_i}{\sum_{i}^{n} x_i^2}$

In this case $\sigma_u^2 = \dfrac{\sum e_i^2}{(n-1)}$ (a) $SE(\beta) = \sqrt{\sigma_u^2 / \sum x_i^2}$ and $R = \dfrac{\sum e_i^2}{\sum y^2}$

## 2.8 Estimation of Plasticities from an Estimated Regression Line

The estimated regression equation is $\hat{Y} = \hat{\alpha} + \hat{\beta}X$ where $\hat{\alpha}$, $\hat{\beta}$ [...] and [...]

Now the derivative of $\hat{Y}$ with respect to $X$ is $\hat{\beta} = \dfrac{d\hat{Y}}{dX}$ which shows the rate of change

in $\hat{Y}$ as $X$ changes by a very small amount. It should be clear that the elasticity [...] function is a linear demand or supply function, the coefficient [...] is not the price elasticity but a component of the elasticity which is defined by the formula

$$\eta_p = \frac{dY}{dX} \cdot \frac{X}{Y} = \frac{X}{Y} \cdot \frac{dY}{dX}$$

where $\eta_p$ = price elasticity, $Y$ = quantity (demanded or supplied), $X$ = price and

$\hat{\beta}$ is the component $\dfrac{dY}{dX}$. From an estimated function we can obtain the average

elasticity $\eta_p = \hat{\beta}\dfrac{\bar{X}}{\bar{Y}}$ where $\bar{X}$ = the average price in the sample, $\bar{Y}$ = average regression

value of the quantity, i.e. the mean value as estimated from the regression, $\bar{Y}$ is [...]

average value of the quantity in the sample. It should be noted that $\bar{Y} = \hat{\bar{Y}}$ since the

$$\hat{Y} = \hat{\alpha} + \hat{\beta}X$$

$$\hat{\bar{Y}} = \hat{\alpha} + \hat{\beta}\bar{X} = \bar{Y} \quad (\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X} = \bar{Y})$$

In particular if $Y = \alpha - \beta X$, is the regression equation then the estimated average

elasticity $\eta_p = \hat{\beta}\dfrac{\bar{X}}{\bar{Y}}$ where $\bar{Y} = \hat{\alpha} - \hat{\beta}\bar{X}$

Now substituting for $\bar{Y}$ in the expression of elasticities we obtain $\eta_p = \dfrac{\beta X}{\alpha + \beta X}$

If the function $Y = \alpha - \beta X$ represents a supply function with $\beta > 0$, it follows that

(i) the supply function will be elastic ($\eta_p > 1$) if $\alpha$ is negative $\alpha < 0$

(ii) the supply function will be inelastic ($\eta_p < 1$) if $\alpha > 0$

(iii) the supply function will have unitary elasticity $\eta_p = 1$ if $\alpha = 0$

Thus the elasticity of a supply curve (with positive slope) depends on the sign of
the constant intercept $\alpha$

**Example 2.2** The following table includes the price and quantity demanded of the
product of a monopolist over a six year period.

| Year | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 |
|------|------|------|------|------|------|------|
| Quantity (000 Kg.) | 5 | 3 | 4 | 1 | 8 | 0 |
| Price (00 ₹) | 2 | 4 | 3 | 1 | 3 | 5 |

(a) Estimate the demand function, assuming a linear demand function. Comment on
   the values of the estimated coefficients ($\alpha$ and $\beta$) on the basis of economic
   theory

(b) Estimate the average elasticity of demand

(c) Estimate the elasticity of demand at the price 4

(d) Forecast the level of demand if price rises to 5 Comment on your forecast

**Solution** (a) Let $Y = \alpha + \beta X$ for $i = 1, 2, \dots 6$ be the linear demand function. By the OLS method we can get the estimators of $\alpha$ and $\beta$. Here $Y$ = demand, $X$ = price $\alpha$, $\beta$ are two parameters. Theoretically we may assume $\alpha > 0$, $\beta < 0$. By OLS method

$$\hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2} \quad \text{where} \quad x_i = X_i - \bar{X}, \quad y_i = Y_i - \bar{Y} \quad \text{and} \quad \alpha = \bar{Y} - \beta \bar{X}, \quad \bar{X} = \frac{\sum X_i}{n}$$

$$\bar{Y} = \frac{\sum Y_i}{n}$$

### Calculations for the parameters $(\alpha, \beta)$

| Year (n) | quantity ( '000 kg ) $Y_i$ | price ( '00 ₹) $X_i$ | $y_i = Y_i - \bar{Y}$ | $x_i = X_i - \bar{X}$ | $x_i y$ | $x^2$ |
|---|---|---|---|---|---|---|
| 20 4(1) | 8 | 2 | 3 | 1 | 3 | 1 |
| 2015 (2) | 3 | 4 | 2 | 1 | 2 | 1 |
| 20 6 (3) | 4 | 3 | -1 | 0 | 0 | 0 |
| 20 7 (4) | 7 | 1 | -2 | 2 | 4 | 4 |
| 20 8 (5) | 8 | 3 | 3 | 0 | 0 | 0 |
| 2019 (6) | 0 | 5 | -4 | 2 | 10 | 4 |
| $n = 6$ Total | $\sum Y$ 30 | $\sum X = 18$ | $\sum y_i = 0$ | $\sum x_i$ 0 | $\sum x_i y = 19$ | $\sum x^2 = 0$ |

$$\bar{Y} = \frac{\sum Y_i}{n} = \frac{30}{6} = 5, \quad \bar{X} = \frac{\sum X_i}{n} = \frac{6}{6} = 3$$

Now $\hat{\beta} = \dfrac{\sum x_i y}{\sum x_i^2} = \dfrac{19}{10} = 1.9$

and $\hat{\alpha} = \bar{Y} - \beta \bar{X} = 5 - (1.9) \times 3 = 5 + 5.7 = 10.7$

Thus the OLS estimators of $\alpha$ and $\beta$ are $\alpha = 10.7 > 0$ and $\beta = -1.9 < 0$.

Therefore the estimated demand function is $Y = \alpha + \beta X$ or $Y = \bar{Y} = 10.7 - 1.9X$

This is consistent with the theory where we assume $\alpha > 0$ and $\beta < 0$. This clearly shows that there exists an inverse relation between price and demand i.e. the law of demand holds true

(b) The average elasticity (price elasticity of demand) is given by.

$$\eta_p = \beta \frac{\bar{X}}{\bar{Y}} = 1.9 \times \frac{3}{5} = 1.14 \quad \text{or} \quad \eta_p = 1.14 > 1$$

This means that the demand function shows an elastic demand

(c) We have to estimate $\eta_p$ (price elasticity of demand) from the estimated relation

$\hat{Y} = 10.7 - 1.9X$ when price $X = 4$

If $X = 4$, $Y = 10.7 - 1.9 \times 4 = 10.7 - 7.6 = 3.1$

Now $\eta_p$ (at $X = 4$) $= \dfrac{X}{Y} \dfrac{dY}{dX} = \dfrac{4}{3.1} \times 1.9 = 2.45$

$\eta_p$ at $X = 4$ is $2.45 > 1$

This implies that the demand is elastic demand

... the ... using the ... of demand when price rises ... e ... where ...

... estimate demand ...

a) ...                              b) ...

where ...

... t ...

This means has ... of ... from ... to ... demand decreases ... as

This ... shows that as price ... demand decreases

**Example 2.3** The following table shows ten pairs of observations in a particular quantity supplied.

| No. of observations | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| Quantity ... tons | ... | ... | ... | 56 | 57 | 55 | 58 | 67 | 62 | 64 |
| Price (in Rs.) | 9 | ... | 6 | 10 | 9 | 10 | 7 | 8 | 2 | 8 |

a) Assuming a linear supply function estimate the supply function comment on the values of the estimated coefficients ($\alpha$ and $\beta$) on the basis of economic theory.

b) Estimate the average price elasticity of supply.

c) Calculate the elasticity of supply at the price 8

d) Forecast the level of supply when price is 8

**Solution** Let $Y = \alpha + \beta X$ for ... $i = 1, 2, ..., 12$ be the linear supply function. By ... method we can get the estimators of $\alpha$ and $\beta$. Here $Y$ = supply $X$ = price ... are the parameters. Theoretically we may assume $\alpha < 0$ and $\beta > 0$. By OLS method we may get

$$\beta = \frac{\sum z_i}{\sum z^2} \quad \text{where} \quad z_i = Y_i - \bar{Y} \quad ...$$

and $\alpha = \bar{Y} - ... \quad Y = \frac{1}{n}\sum Y, \quad \bar{Y} = \frac{1}{n}\sum ...$

**Calculations for the OLS estimators of parameters ($\alpha$, $\beta$)**

| Observations n | $Y$ Quantity (in tons) | $X_i$ Price (in 00 Rs.) | $z_i = ...$ | $Y, \bar{Y}$ | $z_i$ | |
|---|---|---|---|---|---|---|
| 1 | 44 | 9 | 0 | 6 | 0 | |
| 2 | 70 | 12 | 3 | 13 | 0 | |
| 3 | 42 | 6 | 3 | 11 | 33 | 4 |
| 4 | 46 | 10 | | 7 | | |
| 5 | 57 | 9 | 0 | 6 | | |
| $n = 12$ 6 | 77 | 10 | 1 | 14 | 14 | |
| 7 | 58 | 7 | 2 | 5 | 10 | 4 |
| 8 | 45 | 8 | 1 | 8 | 8 | |
| 9 | 67 | 2 | 3 | 4 | 2 | |
| 10 | 57 | 6 | 1 | 10 | 40 | 4 |
| 11 | 72 | 1 | 7 | 9 | 18 | 4 |
| 12 | 64 | 8 | 1 | 1 | | |
| **Total** | $\sum Y = 756$ | $\sum X_i = ...$ | $\sum z = 0$ | $\sum ... = 0$ | $\sum ... = 56$ | $\bar{Y} = ...$ |

$$\bar{x} = \frac{\sum x}{n}, \qquad \ldots = 0, \qquad \bar{y} = \frac{\sum y}{n} \qquad \ldots$$

ii Now the least square estimators of the regression parameters $\alpha$ and $\beta$ are given by

$$\beta = \frac{\sum \ldots}{\sum x_i} = \frac{54}{44} \quad 1.25$$

and $\gamma = \bar{y} - \bar{x}$ ... $\quad 1.25 \quad ... \quad ... \quad 33.75$

Thus the estimated supply function is $\bar{y} = \alpha + \beta x$ or $\hat{y} = 33.75 + 1.25x$

Here we see that $\alpha = 33.75 > 0$ and $\beta = 1.25 > 0$. This means that here is a direct 'positive' relation between supply and price. The intercept of the supply function is positive here. Hence our results are consistent with the theory.

b Average price elasticity of supply is given by $\eta_p = \beta \frac{\bar{x}}{\bar{y}} = 1.25 \quad ... \quad .46$

This shows that at the average price the supply is price inelastic.
c We have to find price elasticity of supply at price 6.
Since the estimated supply function is

$y = \alpha + \beta x$ or $y = 33.75 + 1.25x$
Now $E(y = 6) = 33.75 + 1.25 \times 6 = 33.75 \quad ... \quad 53.25$

Now by definition price elasticity of supply $\eta_p = \frac{x}{y} \cdot \frac{dy}{dx} = \frac{6}{53.25} \quad 1.25 \quad ...$

Thus $\eta_p = 0.366$ when $X = 6$.
e From the estimated supply function we see that $y = 33.75 + 1.25x$
When $X = 6$, $Y = 53.25$
If now price increases to 8 i.e., if $X = 8$
then $Y = 33.75 - 3.25 + 8 = 33.75 + 2h = 59.75$
This means that when $X = 6$, $Y = 53.25$
and when $X = 8$, $Y = 59.75$
Thus we may forecast that as price increases supply will also increase.

## 2 9. Properties of Least Squares Estimators

The least squares estimates are called BLUE 'best linear unbiased estimates' provided that the random term $u$ satisfies some general assumptions namely that the $u$ has zero mean and constant variance. This proposition together with the set of conditions under which it is true is known as **Gauss Markov Least-Squares Theorem**.

The OLS estimators possess three properties. They are linear unbiased and have the smallest variance compared to other linear unbiased estimators. Thus the OLS estimators are BLUE.

### 1 The property of linearity

The least-squares estimates $\alpha$ and $\beta$ are linear functions of the observed sample values $Y$.

I'm sorry, but this page is too degraded for me to transcribe reliably.

## 2 The property of unbiasedness

The means of $\hat{\alpha}$ $[E(\hat{\alpha})]$ and $\hat{\beta}$ $[E(\hat{\beta})]$ can be obtained as follows

$$\text{Since } \hat{\beta} = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$$

$$= \frac{\sum_{i=1}^{n} Y_i x_i + \bar{Y}\sum_{i=1}^{n}(x_i - \bar{X})}{\sum_{i=1}^{n}(X_i - \bar{X})^2} = \frac{\sum_{i=1}^{n} x_i Y_i}{\sum_{i=1}^{n} x_i^2} \text{ where } \sum_{i=1}^{n} x_i = (X_i - \bar{X}) = 0$$

and $x_i = X_i - \bar{X}$ for $i = 1, 2, \ldots, n$

$$\hat{\beta} = \sum_{i=1}^{n} K_i Y_i \text{ where } K_i = \frac{x_i}{\sum_{i=1}^{n} x_i^2} \text{ and } \alpha = \sum_{i=1}^{n} \frac{1}{n} (\bar{X})_i$$

Now $\hat{\beta} = \sum_{i=1}^{n} K_i Y_i$ We now put $Y_i = \alpha + \beta X_i + u_i$

$$\hat{\beta} = \sum_{i=1}^{n} K_i (\alpha + \beta X_i + u_i) = \alpha \sum_{i=1}^{n} K_i + \beta \sum_{i=1}^{n} K_i X_i - \sum_{i=1}^{n} K_i u_i$$

Since $K_i = \frac{x_i}{\sum_{i=1}^{n} x_i^2}$ $\sum_{i=1}^{n} K_i = \frac{\sum_{i=1}^{n} x_i}{\sum_{i=1}^{n} x_i^2} = 0,$ as $\sum_{i=1}^{n} x_i = 0$

and $\sum_{i=1}^{n} K_i X_i = \sum_{i=1}^{n} K_i (x_i + \bar{X})$ where $x_i = X_i - \bar{X}$  $X_i = x_i + \bar{X}$

$$= \sum_{i=1}^{n} K_i x_i + \bar{X} \sum_{i=1}^{n} K_i = \sum_{i=1}^{n} K_i x_i \left[ \because \sum_{i=1}^{n} K_i = 0 \right]$$

Now $\sum_{i=1}^{n} K_i x_i = \sum_{i=1}^{n} x_i x_i \Big/ \sum_{i=1}^{n} x_i^2 = \frac{\sum_{i=1}^{n} x_i^2}{\sum_{i=1}^{n} x_i^2} = 1$ as $K_i = \frac{x_i}{\sum_{i=1}^{n} x_i^2}$

1. we put $\sum_i k_i \quad$ and $\sum_i k_i x_i \quad$ in the expression of

$\beta = \sum_i k_i + \sum_i k_i x_i + \sum_i k_i u_i \quad$ we get $\beta = \beta + \sum_i k_i u_i$

Now mean of $\beta = E(\beta) = E(\beta) + \sum_i k_i E(u_i) \quad \beta \quad E(u_i) = 0$

thus we have the $\beta$ i.e. mean of $\beta$ is $\beta$

Similarly $= \sum_n \quad (k_i) \beta$

$\sum_n \quad (k_i) \beta + \alpha + \beta x_i + u_i \quad \beta x_i \quad u_i$

$\sum_n \frac{1}{n} \alpha + \beta \cdot \frac{1}{n} \sum_1 x_i + \frac{1}{n} \sum_{-1} u_i, \quad \text{as} \sum_i k_i \quad 0 \quad \sum_i k_i x_i \quad \sum_i k_i u_i$

Since $\sum_i k_i = 0, \quad \sum_i k_i x_i = 1, \quad \sum_i 1 = n$ we have

$= \alpha + \beta \bar{x} + \frac{1}{n} \sum_i u_i - \beta \bar{x} + \frac{1}{n} \sum_i k_i u_i$

or $\alpha = \alpha + \frac{1}{n} \sum_{i=n} u_i - \bar{x} \sum_i k_i u_i, \quad$ or $E(\alpha) = E(\alpha) + \frac{1}{n} \sum_i E(u_i) - \sum_i k_i E(u_i),$

$E(u_i) = \alpha \quad \text{as} \quad E(u_i) = 0$

This shows that mean of $\alpha$ is $\alpha$

Thus it is proved that $\alpha$ and $\beta$ are unbiased estimators of $\alpha$ and $\beta$.

### 1 The minimum variance property

In this property we shall prove the Gauss Markov Theorem which states that the least squares estimates are best have the smallest variance as compared with any other linear unbiased estimator obtained from other econometric methods.

First we have to find var $(\beta)$ and var $(\alpha)$ and then we have to prove the minimum variance property.

Variance of $\beta = \text{var}(\beta) = E[\beta - E(\beta)]^2 = E[\beta - \beta]^2$ as $E(\beta) = \beta$

Since $\hat{\beta} = \beta + \sum_{i}^{n} K_i u_i$ (see property 2)

$$\hat{\beta} - \beta = \sum_{i}^{n} K_i u_i$$

or $(\hat{\beta} - \beta)^2 = \left(\sum_{i}^{n} K_i u_i\right)^2$   or $E(\hat{\beta} - \beta)^2 = E\left(\sum_{i}^{n} K_i u_i\right)^2$

or $var(\hat{\beta}) = E\left[\sum_{i}^{n} K_i^2 u_i^2 + 2\sum_{i \neq j}\sum K_i K_j u_i u_j\right]$

$$= \sum_{i}^{n} K_i^2 E(u_i^2) + 2\sum_{i \neq j}\sum K_i K_j E(u_i u_j)$$

$$= \sum_{i=1}^{n} K_i^2 E(u_i^2) \qquad [\because E(u_i u_j) = 0, \text{ for } i \neq j$$

$$= \sum_{i=1}^{n} K_i^2 \sigma_u^2 \qquad [\because E(u_i^2) = \sigma_u^2$$

$$= \frac{\sum_{i}^{n} x_i^2}{\left(\sum_{i}^{n} x_i^2\right)^2} \sigma_u^2 \qquad \left[\because K_i = \frac{x_i}{\sum x_i^2}\right]$$

$$= \frac{\sigma_u^2}{\sum x_i^2} \qquad \therefore var(\hat{\beta}) = \frac{\sigma_u^2}{\sum x_i^2}$$

**Similarly, Variance of $\alpha = var \hat{\alpha}$**

Since $\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$ (see property 1)

Substituting $\hat{\beta} = \sum_{i=1}^{n} K_i Y_i$ we obtain $\hat{\alpha} = \bar{Y} - \bar{X}\sum_{i=1}^{n} K_i Y_i$

$$= \frac{\sum_{i}^{n} Y_i}{n} - \bar{X}\sum_{i=1}^{n} K_i Y_i = \sum_{i=1}^{n}\left(\frac{1}{n} - \bar{X}K_i\right)Y_i$$

$$\text{Now} \quad var\ \alpha = var \sum \frac{1}{n}\ (A_i) = \sum_{i=1}^{n}\left[\frac{1}{n}\ (A_i)\right]\ var\ (\ )$$

Since $var\ (\ ) = \sigma_u$

$$var\ \alpha = \sum_{i}^{n}\frac{1}{n}\ (A)\ \sigma_u^2$$

$$\sigma_u \sum_{i}^{n}\frac{1}{n}\ (\bar{X}A) = \sigma_u \sum_{i=1}^{n}\frac{1}{n} : \frac{1}{n}\ \bar{X}A + \bar{X}^2 A^2$$

$$= \sigma_u\left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^{n}x_i^2}\right] \qquad \sum_{i=1}^{n}X_i = 0 \text{ and } \sum_{i=1}^{n}A^2 = \frac{1}{\sum_{i=1}^{n}x_i^2} \text{ where } A = \frac{x}{\sum_{i=1}^{n}x_i}$$

$$= \sigma_u^2\left[\frac{\sum_{i=1}^{n}x_i^2 - n\bar{x}^2}{n\sum_{i=1}^{n}x_i^2}\right] \cdot \sigma_u^2 \frac{\sum_{i=1}^{n}x_i^2}{n\sum_{i=1}^{n}x_i^2}$$

$$For \quad \text{i)} \quad \sigma_u \sum_{i}^{n}x_i^2 \bigg/ n\sum_{i=1}^{n}x_i^2$$

$$\left[\sum_{i=1}^{n}x_i^2 + n\bar{x}^2 = \sum_{i=1}^{n}(X - \bar{X})^2 + n\bar{x}^2 = \sum_{i=1}^{n}x_i^2 - 2\bar{X}\sum_{i=1}^{n}X + n\bar{x}^2 + n\bar{x}^2\right.$$

$$\left. \sim \sum_{i=1}^{n}X_i - 2n\bar{x}^2 + 2n\bar{x}^2 = \sum_{i=1}^{n}X_i^2\right]$$

**Case (a) $\beta$ has the least variance .**

We know that $Var(\beta) = \sigma_u^2 \bigg/ \sum_{i=1}^{n}x_i^2$

Now we want to prove that any other linear unbiased estimate of the true parameter for example $\beta^*$ obtained from any other econometric method, has a bigger variance than the least squares estimate $\beta$. Thus we have to prove that    $var\ \beta < var\ \beta^*$

**Proof** The new estimator $\beta^*$ is by assumption a linear combination of the Y, a

weighted sum of the sample values $Y_i$, the weights $K^* = x \sum_{i}^{n}x_i^2$ being different from the weights of the least-squares estimates

For example, let us assume $\beta^* = \sum c_i Y_i$, where $c_i = k_i + d_i$, $d_i$ is an arbitrary set

weights sum as (but not the same) to the $k_i$'s

$c_i$ is put $Y_i = \alpha + \beta X_i + u_i$, in the expression of $\beta^*$ and we obtain

$$\beta^* = \sum_{i=1}^{n} c_i(\alpha + \beta X_i + u_i) = \sum (\alpha c_i + \beta c_i X_i + c_i u_i)$$

It is assumed that the $\beta$, $\beta^*$ is also an unbiased estimator of $\beta$, i.e. $E(\beta^*) = \beta$

Now $E(\beta^*) = E\left[\sum_{i=1}^{n} \alpha c_i + \beta c_i X_i + c_i u_i\right]$

$E(\beta^*) = E\left[\alpha \sum_{i=1}^{n} c_i + \beta \sum_{i=1}^{n} c_i X_i + \sum_{i=1}^{n} c_i u_i\right]$

Now $E(\beta^*) = \beta$ if, and only if

$$\sum_{i=1}^{n} c_i = 0, \quad \sum_{i=1}^{n} c_i X_i = 1 \text{ and } \sum_{i=1}^{n} c_i u_i = 0$$

But $\sum_{i=1}^{n} c_i = 0$ implies $\sum_{i=1}^{n} d_i = 0$ because

$$\sum c_i = \sum (K_i - d_i) = \sum_{i=1}^{n} K_i + \sum_{i=1}^{n} d_i, \text{ and } \sum_{i=1}^{n} K_i = \frac{\sum_{i=1}^{n} x_i}{\sum_{i=1}^{n} x^2} = 0 \text{ as } \sum x_i = 0$$

i.e. $\sum_{i=1}^{n} c_i = \sum_{i=1}^{n} d_i$. Therefore if $\sum_{i=1}^{n} c_i = 0$, then $\sum_{i=1}^{n} d_i = 0$

Similarly $\sum_{i=1}^{n} c_i X_i = 1$ requires $\sum_{i=1}^{n} d_i X_i = 0$,

since $\sum_{i=1}^{n} c_i X_i = \sum_{i=1}^{n} (K_i + d_i) X_i = \sum_{i=1}^{n} K_i X_i + \sum_{i=1}^{n} d_i X_i$

Given that $\sum_{i=1}^{n} K_i X_i = 1$, $\sum_{i=1}^{n} c_i X_i = 1$ if $\sum_{i=1}^{n} d_i X_i = 0$

Thus $\beta^*$ will be a linear unbiased estimate of $\beta$ (with weights $c_i = K_i + d_i$, if

$$\sum_{i=1}^{n} c_i = 0, \quad \sum_{i=1}^{n} d_i = 0, \quad \sum_{i=1}^{n} c_i X_i = 1 \text{ and } \sum_{i=1}^{n} d_i X_i = 0$$

Note ...

$$\hat{\beta} = \sum k_i \ldots \qquad \sum k_i$$

$$\sum_i a_i k_i \quad \sum_i c_i \quad \sum_i k_i \ldots$$

for $\hat{\beta} = \beta + \sum \ldots$ where $\ldots = \beta(1 + \ldots$.

$E(\ldots) = E \ldots = \beta(1 + \ldots)$

$= \beta(1 \ldots)$

Now ...

$$\ldots = \sigma_u^2$$

Similarly we may obtain

$$\beta^* = \sum \ldots \quad \text{and} \quad Var(\beta^*) = Var\left(\sum \ldots - \sum \ldots \right) \quad Var\, \beta = \sum \ldots \sigma_u^2$$

Now $\sum_{i=1} \ldots \sum (k_i + d_i)$

$$\sum_{i=1} k_i + \sum d_i + 2\sum k_i d_i = \sum k_i \sum d_i^2$$

Given that $\sum k_i d_i = \dfrac{\sum x_i d_i}{\sum x_i^2} = \dfrac{\sum (X - \bar{X}) d_i}{\sum x_i^2}$

$$= \dfrac{\sum d_i X - \bar{X}\sum d_i}{\sum x_i^2} = 0 \quad \left(\text{as } \sum d_i X = 0 \text{ and } \sum d_i = 0\right)$$

Substituting we find

$$Var(\beta^*) = \sigma_u^2 \left( \sum k_i + \sum d_i^2 \right) = \sigma_u^2 \sum k_i^2 + \sigma_u^2 \sum d_i^2$$

$$Var(\beta^*) = Var(\beta) + \sigma_u^2 \sum_i d_i^2 \qquad Var(\beta) = \sigma_u^2 \sum k_i^2 = \sigma_u^2 \sum x_i^2$$

Since $\sigma_u \sum$ ... it proves that $Var(\beta^*) \geq Var(\hat\beta)$ or $Var \geq Var \beta^*)$

Thus it is proved that $\hat\beta$ is the BLUE of $\beta$.

**Case (b)** In the same way it can be proved that the least squares constant intercept $\alpha$ possesses minimum variance. We take a new estimator $\alpha^*$ which we assume to be a linear function of the $Y$'s and we obtain $c_i = A_i + d_i$

where $A_i$
$$\sum c_i$$

Since $\alpha \quad \sum_n \quad (A_i)$

Similarly, $\alpha^* = \sum_{i=1}^{n} \left(\frac{1}{n} \quad \bar X c_i\right) Y_i$. (1)

This shows that like $\alpha$, $\alpha^*$ is also a linear function in $Y$'s
Now $\alpha^*$ is to be regarded as an unbiased estimator of $\alpha$ ( $E(\alpha^*) = \alpha$
We substitute for $Y_i = \alpha + \beta X_i + u_i$ in $\alpha^*$ and we get,

$$\alpha^* = \alpha \left[ \bar X \sum_i c_i \right] - \beta \left[ X - \bar X \sum_{i=1}^{n} c_i X_i \right] - \sum_{i=1}^{n}\left[\frac{1}{n} - \bar X c_i\right] u_i$$

Now $E(\alpha^*) = \alpha \left[1 - \bar X \sum_{i=1}^{n} c_i\right] + \beta\left[\bar X - \bar X \sum_{i=1}^{n} c_i X_i\right] + E\left[\sum_{i=1}^{n}\frac{1}{n} - \bar X c_i\right] u_i$

Now $E(\alpha^*) = \alpha$ if and only if $\sum_{i=1}^{n} c_i = 0$, $\sum_{i=1}^{n} c_i X_i = 1$ and $\sum_{i=1}^{n} c_i u_i = 0$

These conditions imply $\sum_{i=1}^{n} d_i = 0$ and $\sum_{i=1}^{n} d_i X_i = 0$

The variance of $\alpha^*$ is given by
$Var(\alpha^*) = E[\alpha^* - E(\alpha^*)]^2 = E[\alpha^* - \alpha]^2$

$$\sigma_u^2 \sum_{i=1}^{n}\left[\frac{1}{n} - \bar X c_i\right]^2 = \sigma_u^2 \sum_{i=1}^{n}\left[\frac{1}{n^2} - \frac{2}{n}\bar X c_i + \bar X^2 c_i^2\right]$$

$$= \sigma_u^2\left[\frac{n}{n^2} - 2\bar X \frac{1}{n}\sum_{i=1}^{n} c_i + \bar X^2 \sum_{i=1}^{n} c_i^2\right] = \sigma_u^2\left[\frac{1}{n} + \bar X^2 \sum_{i=1}^{n} c_i^2 - \frac{2}{n}\bar X \sum_{i=1}^{n} c_i\right]$$

Since $\sum_{i=1}^{n} c_i = 0$ and $\sum_{i=1}^{n} c_i^2 \quad \sum_{i=1}^{n} A_i^2 + \sum_{i=1}^{n} d_i^2$

we have $\text{var } a^* = \sigma_u \left[ \frac{1}{n} + \bar{X}^2 \sum A + \sum d \right]$

$$= \sigma_u^2 \left[ \frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] + \sigma_u^2 \left[ \sum d_i \right], \text{ where } \sum A = \sum x_i^2$$

$$\text{var}(a^*) = \sigma_u^2 \left[ \frac{1}{n} + \frac{\bar{X}^2}{\sum x_i^2} \right] + \sigma_u^2 \bar{X}^2 \sum d_i^2, \text{ or, } \text{var}(a^*) = \text{var}(\hat{a}) + \sigma_u^2 \bar{X}^2 \sum d^2$$

Here $\sum d^2 > 0$, because all $d_i$ s are not zero

Thus we have, $\text{var}(a^*) > \text{var}(a)$ or, $\text{var}(a) < \text{var}(a^*)$

Hence it is proved that $\hat{a}$ is the BLUE of $\alpha$.

## 2.10. The Variance of the Random Variable, u

The formulae of the variance of $\hat{a}$ and $\beta$ involve the variance of the random term u, $\sigma_u^2$. However, the true variance of $u$, cannot be computed since the values of $u_i$ are not observable. But we may obtain an unbiased estimate of $\sigma_u^2$ from the expression

$$\hat{\sigma}_u^2 = \sum_{i=0}^{n} e_i^2 \bigg/ (n-2) \text{ where } e_i = Y_i - \hat{Y}_i = Y_i - \hat{a} - \beta X_i,$$

[$Y$ is the observed value and $\hat{Y}_i$ is the estimated value i.e. $Y = \alpha + \beta X_i + e_i$, and $\hat{Y} = \hat{a} + \beta X_i$ for $i = 1, 2, \ldots, n$]

**Proof** One property of the regression line $\hat{Y}_i = \hat{a} + \beta X_i$, is that it passes through the point $(\bar{X}, \bar{Y})$. So, $\bar{Y} = \hat{a} + \beta \bar{X}$

Again we know that $\bar{Y} = \alpha + \beta \bar{X} + \bar{u}$ from the observed relationship

$$\left[ \text{Where } Y_i = \alpha + \beta X_i + u_i \quad \sum_{i=1}^{n} Y_i = n\alpha + \beta \sum_{i=1}^{n} X_i + \sum_{i=1}^{n} u_i \right.$$

$$\left. \text{or, } \sum_{i=1}^{n} Y_i \bigg/ n = \alpha + \beta \sum_{i=1}^{n} X_i \bigg/ n + \sum_{i=1}^{n} u_i \bigg/ n \text{ or, } \bar{Y} = \alpha + \beta \bar{X} + \bar{u} \right]$$

Since $e_i = Y_i - \hat{Y}_i$

$= Y_i - [\hat{\alpha} + \hat{\beta} X_i] \quad (\alpha + \beta X_i + u_i - \hat{\alpha} - \hat{\beta} X_i \quad \bar{u} \quad \alpha - \beta \bar{X} \quad \alpha \quad \beta \bar{X})$

$= [\beta(X_i - \bar{X}) + (u_i - \bar{u})] - [\hat{\beta}(X_i - \bar{X})]$

$= -[(\hat{\beta} - \beta)x_i + (u_i - \bar{u})] \quad$ where $x_i = X_i - \bar{X}$

or $e_i^2 = (\hat{\beta} - \beta)^2 x_i^2 + (u_i - \bar{u})^2 - 2x_i(\hat{\beta} - \beta)(u_i - \bar{u})$

$= (\hat{\beta} - \beta)^2 x_i^2 + (u_i - \bar{u})^2 - 2(\hat{\beta} - \beta)x_i(u_i - \bar{u})$

$\sum_{i=1}^{n} e_i^2 = (\hat{\beta} - \beta)^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n}(u_i - \bar{u})^2 - 2(\hat{\beta} - \beta)\sum_{i=1}^{n} x_i u_i - n\sum_{i=1}^{n} e_i$

$$= \frac{\sum_{i=1}^{n} x_i u_i}{\sum_{i=1}^{n} x_i^2} \quad \sum_{i=1}^{n} x_i^2 + \left[\sum_{i=1}^{n} u_i^2 - \frac{\left(\sum_{i=1}^{n} u_i\right)^2}{n}\right] - 2\frac{\sum_{i=1}^{n} x_i u_i}{\sum_{i=1}^{n} x_i^2}$$

Since $\hat{\beta} = \beta + \dfrac{\sum_{i=1}^{n} x_i u_i}{\sum_{i=1}^{n} x_i^2} \qquad \hat{\beta} - \beta = \sum_{i=1}^{n} x_i u_i \Big/ \sum_{i=1}^{n} x_i^2$

Again, $\sum_{i=1}^{n}(u_i - \bar{u})^2 = \sum_{i=1}^{n} u_i^2 - 2\bar{u}\sum_{i=1}^{n} u_i + \sum_{i=1}^{n} \bar{u}^2$

$= \sum_{i=1}^{n} u_i^2 - 2\bar{u} \cdot n \cdot \frac{1}{n}\sum_{i=1}^{n} u_i + n\bar{u}^2 = \sum_{i=1}^{n} u_i^2 - 2n\bar{u}^2 + n\bar{u}^2$

$= \sum_{i=1}^{n} u_i^2 - n\bar{u}^2 = \sum_{i=1}^{n} u_i^2 - n\left(\frac{\sum_{i=1}^{n} u_i}{n}\right)^2 = \sum_{i=1}^{n} u_i^2 - \frac{\left(\sum_{i=1}^{n} u_i\right)^2}{n}$

and $2(\hat{\beta} - \beta)\sum_{i=1}^{n} x_i u_i - \bar{u}\sum_{i=1}^{n} x_i = 2\dfrac{\sum_{i=1}^{n} x_i u_i}{\sum_{i=1}^{n} x_i^2}\sum_{i=1}^{n} x_i u_i = 2\dfrac{\left(\sum_{i=1}^{n} x_i u_i\right)^2}{\sum_{i=1}^{n} x_i^2}$ as $\bar{u} = 0$

$$\sum_{i=1}^{n} e_i^2 = \left\{\sum_{i=1}^{n} u_i^2 - \frac{\left(\sum_{i=1}^{n} u_i\right)^2}{n}\right\} - \left(\sum_{i=1}^{n} x_i u_i\right)^2 \Big/ \sum_{i=1}^{n} x_i^2$$

$$\text{of } \sum_{i}^{n} \qquad \sum_{i}^{n}$$

$$E\left[\sum_{j}^{n} \qquad \sum_{j}^{n} \right]$$

$$= \left[\sum_{i} \sigma_u^2 \quad \sum_{i}^{n} \frac{\sigma_u}{N} \quad \sum_{i} \quad \sum_{i} \quad \left( \qquad \right) \right]$$

$$\text{or } E\left[\sum_{i=1}^{n} e^2 \right] = N\sigma_e^2 \quad \frac{N\sigma_u^2}{n} \quad \sigma_u \sum_{i} \quad \sum_{i}$$

$$\left( \sigma_u^2 \quad \sigma_v^2 \quad \sigma_v^2 \right)$$

or, $E\left[\dfrac{\sum_{i} e_i^2}{n-2}\right] = \sigma_u^2$

So, $\sum_{i} e_i^2 / (n-2)$ is an unbiased estimator of $\sigma_u$. If we define

$\sum_{i} e_i^2 / (n-2) = \sigma_u^2$ then $\sigma_u^2$ is an unbiased estimator of $\sigma_u^2$.

## 2.8. Maximum Likelihood Estimators (MLE's) of α, β and $\sigma_u^2$

If each $u_i | X_i = \alpha + \beta X_i + u_i$ is normally distributed with mean 0 and variance $\sigma_u$ i.e., $u_i \sim N(0, \sigma_u^2)$ and all $u_i$ are independent, then MLE of α and β are equal to the OLS estimators of α and β (i.e. $\hat\alpha$ and $\hat\beta$).

**Proof** Since $u_i \sim N(0, \sigma_u^2)$, the p.d.f. of $u_i$ is given by

$$f_i(u_i) = \frac{1}{\sqrt{2\pi}\sigma_u} e^{-\frac{1}{2}\left(\frac{u_i}{\sigma_u}\right)^2} \qquad \frac{1}{\sqrt{2\pi}\sigma_u} e^{\frac{1}{2}x} \qquad \text{as } u_i = 0.$$

This joint probability distribution function of $u_1, u_2, \dots, u_n$ is given by $f(u_1, \dots, u_n)$ and given the set of sample observations it is looked upon as a function of b

parameters and is called the ... form of the parameters ... if ... $x_1 ... x_n$ are independent, then ...

$$f(u_1 ... u_n) = ... u_n) = ... b ...$$

$$w f(u_1 ... \sigma_u) ... \frac{1}{...} f(u ... \beta, \sigma_u) ... $$

of $... \beta, \sigma_u) = \frac{1}{... \sigma_u} ...$

Taking Log on both sides we get,

$$\log f \quad \frac{n}{2} \log 2\pi \quad n \log \sigma_u \quad \frac{1}{2\sigma_u} \sum u ...$$

$$= \frac{n}{2} \log 2\pi \quad n \log \sigma_u \quad \frac{1}{2\sigma_u} \sum_{i=1}^{n} (Y_i \quad \alpha \quad \beta X_i)$$

[ ] $= \alpha - \beta X_i + u_i$ or $u_i = Y_i \quad \alpha - \beta X_i$ or $\sum u_i = \sum Y_i \quad n \quad \beta X_i$

M.L.E of $\alpha$ and $\beta$ can be obtained by maximizing $\log L$ through the choice of $\alpha$ and $\beta$. Maximisation of $\log L$ through the choice of $\alpha$ and $\beta$ is equivalent to minimisation

of $\sum_{i=1}^{n} (Y_i - \alpha - \beta X_i)^2$ through the choice of $\alpha$ and $\beta$.

Let us suppose that $\alpha^*$ is the MLE of $\alpha$ and $\beta^*$ is the MLE of $\beta$. Then,

$$\beta^* = \frac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2} \quad \alpha^* = \bar{Y} \quad \beta^* \bar{X} \qquad \left[ \text{Since } \beta = \sum_{i=1}^{n} ... \quad \sum_{i=1}^{n} ... \text{ and } ... = \bar{Y} \quad \beta \bar{X} \right]$$

Since $\log L \quad \frac{n}{2} \log 2\pi \quad n \log \sigma_u \quad \frac{1}{2\sigma_u^2} \sum_{i=1}^{n} (Y_i \quad \alpha \quad \beta X_i)^2$

Differentiating partially $\log L$ with respect to $\alpha$, $\beta$
we get,

$$\frac{\delta \log L}{\delta \alpha} \quad \frac{1}{2\sigma_u^2} 2 \sum_{i=1}^{n} (Y_i \quad \alpha \quad \beta X_i)(-1)$$

$$\frac{\delta \log L}{\delta \log \beta} = \frac{1}{2\sigma_u^2} \sum_{i=1}^{n} (Y_i \quad \alpha \quad \beta X_i)(-X_i)$$

Equating these equations to zero and putting star mark on the parameter to distinguish them from least squares estimators

we get

$$\frac{\partial}{\partial \alpha_v} \sum \cdots = \alpha^* + \beta^* X_i \qquad (1)$$

$$\frac{1}{\partial \sigma_v} \sum (Y_i) - \alpha^* + \beta^* X_i = 0 \qquad (2)$$

The first two equations are reduced to the least squares normal equations

$$\sum_{i=1}^{n} Y_i = n\alpha^* + \beta^* \sum_{i=1}^{n} X_i$$

$$\sum_{i=1}^{n} X_i Y_i = \alpha^* \sum_{i=1}^{n} X_i + \beta^* \sum_{i=1}^{n} X_i$$

Now solving the two normal equations we can get $\beta^* = \dfrac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$ and $\alpha^* = \bar{Y} \cdots$

This proves that the MLE of $\alpha$ and $\beta$ are the same as the least squares estimators. Hence, they would also possess all the desirable properties.

When $\log \ldots$ is maximised through the choice of $\alpha$ and $\beta$ and $\alpha^*$ and $\beta^*$ are the MLE of $\alpha$ and $\beta$, then,

$$\sum_{i=1}^{n} u_i = \sum_{i=1}^{n} Y_i - \alpha^* - \beta^* X_i)^2 \qquad [\because \alpha^* = \alpha \text{ and } \beta^* = \beta]$$

Hence,    $$\sum_{i=1}^{n} u_i^2 = \sum_{i=1}^{n} (Y_i - \alpha - \beta X_i)^2 = \sum_{i=1}^{n} (Y_i - \hat{Y}_i)^2 = \sum_{i=1}^{n} e_i^2 \qquad [\because \hat{Y} = \alpha \cdots]$$

$$\sum_{i=1}^{n} u_i^2 = \sum_{i=1}^{n} e_i^2$$

So, the likelihood function maximised with respect to $\alpha$ and $\beta$ is given by

$$\log L = -\frac{n}{2} \log 2\pi - n \log \sigma_u - \frac{1}{2\sigma_u^2} \sum e_i^2$$

In order to obtain the MLE of $\sigma_u^2$, $\log L$ is to be maximised through the choice of $\sigma_u$ and the first order condition of maximisation is given by (assuming $\sigma_u^{2*}$ as the MLE of $\sigma_u^2$).

$$\frac{\partial \log L}{\partial \sigma_u} = -\frac{n}{\sigma_u} + \frac{1}{2} \sum_{i=1}^{n} e_i^2 (-2) \frac{1}{\sigma_u^3} = 0 = \frac{1}{\sigma_u}\left[ -n + \sum_{i=1}^{n} e_i^2 \Big/ \sigma_u^2 \right] = 0$$

or, $\sum_{i=1}^{n} e_i^2 \Big/ \sigma_u^2 = n$   or,  $\sigma_u^2 = \sum_{i=1}^{n} e_i^2 \Big/ n = \sigma_u^{2*}$ (say)

so $\sum_{i=1}^{n} \hat{e}_i^2 / n = \sigma_u^2 *$ says to be MLE of the variance of the disturbance term

denoted by $\sigma_u^2 *$

Note. For maximisation, however, we require second order conditions which is not shown here. But we should also check the second order conditions required for maximisation. This means that we have to show that $\dfrac{\delta^2 l}{\delta \alpha^2} < 0$, $\dfrac{\delta^2 l}{\delta \beta^2} < 0$ and $\dfrac{\delta^2 l}{\delta \sigma_u^2} < 0$

Thus we see that the MLE of $\sigma_u^2$ i.e. $\sigma_u^2 * \sum_{i=1}^{n} \hat{e}_i^2 / n$ is not an unbiased estimator but it is a consistent estimator.

i.e. $E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / n \right] = E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / (n-2) \right] \cdot \dfrac{n-2}{n}$

$= E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / (n-2) \left( \dfrac{n-2}{n} \right) \right] = E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / (n-2) \left( 1 - \dfrac{2}{n} \right) \right]$

$= \sigma_u^2 \left( 1 - \dfrac{2}{n} \right)$  since $E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / (n-2) \right] = \sigma_u^2$

$E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / n \right] = \sigma_u^2 \left( 1 - \dfrac{2}{n} \right)$

Now $E\left[ \sum_{i=1}^{n} \hat{e}_i^2 / n \right] \to \sigma_u^2$ as $n \to \infty$

This proves that the MLE of $\sigma_u^2$ i.e. $\sum_{i=1}^{n} \hat{e}_i^2 / n$ is a consistent estimator of $\sigma_u^2$

Note : MLE of $\alpha$ and $\beta$ i.e., $\alpha*$ and $\beta*$ are unbiased estimators of $\alpha$ and $\beta$ but MLE of $\sigma_u^2$ i.e. $\sigma_u^2 * = \sum_{i=1}^{n} \hat{e}_i^2 / n$ is not an unbiased estimator rather it is a consistent estimator (consistently unbiased) of $\sigma_u^2$

## 2.12 The Sampling Distribution of the Least Squares Estimates

Since least squares estimators are linear combinations of independent normal variables, $Y_1, Y_2, \ldots Y_n$, $\hat{\alpha}$ and $\hat{\beta}$ must also be normally distributed with the following characteristics

... ... are of ... ... ... ... mean being equal to their true value of $\alpha$ and $\beta$.

... ... ... ...

... ... ... ... ... ...

$$\alpha \sim N\left(\alpha, \; \sigma_u^2 \; \frac{1}{\sum}\right) \qquad var\left(\; \sigma_u^2 \; \frac{1}{\sum}\right)$$

$$\beta \sim N\left(\beta, \; \frac{\sigma_u^2}{\sum}\right) \quad \cdot E\left[\beta\right] = \beta \quad var(\beta) \quad \frac{\sigma_u^2}{\sum}$$

Variances of the parameters are directly related to the variances of the disturbance. The following points should be noted carefully.

(i) Larger the value of $\sigma_u^2$, the larger the variances of $\alpha$ and $\beta$. In other words the greater the dispersion of the disturbance terms around the population regression line, the greater is the dispersion in the values of estimated regression parameters.

(ii) $\sum x^2$ in the denominator of the variance formula of both the estimators. This indicates that the more dispersed the values of the explanatory variables are, we will get $\sum x^2$ the smaller the variances of $\alpha$ and $\beta$. If $\sum x$ tends to zero i.e. when all $x = \bar{x}$ then both variances would be infinitely large.

(iii) the variance of $\alpha$ is the smallest when $\bar{x} = 0$ or tends to zero, in particular when $\bar{x} = 0$; $var(\alpha) = \frac{\sigma_u^2}{N}$

## 2.13 Confidence Intervals and Hypothesis Testing

It is highly essential to construct confidence intervals of the parameters in order to achieve precision of $\alpha$ and $\beta$. We have all the information concerning the distribution of $\alpha$ and $\beta$ in order to standardise them.

Since $\alpha \sim N\left(\alpha, \sigma_u^2\left[\frac{1}{N} + \frac{\bar{x}^2}{\sum\limits x^2}\right]\right)$ and $\beta \sim N\left(\beta, \frac{\sigma_u^2}{\sum\limits_{i=1} x^2}\right)$

$$\text{var}(\hat{\beta}_0) = \sigma_u^2 \left[ \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^{n} x_i^2} \right]$$

where $SE(\cdot) = \sqrt{\text{var}(\cdot)}$

Here $\text{var}(\hat{\beta}_0) = \sigma_u^2 \left[ \frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^{n} x_i^2} \right]$

$$SE(\hat{\beta}_0) = \sqrt{\text{var}(\hat{\beta}_0)} = \sigma_u \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^{n} x_i^2}} = \sigma_u \sqrt{\sum_{i=1}^{n} X_i^2 \Big/ n\sum_{i=1}^{n} x_i^2}$$

$$\frac{n + \frac{\bar{X}^2}{\sum_{i=1}^{n} x_i^2}} = \frac{\sum_{i=1}^{n} x_i^2 - n\bar{X}^2}{n\sum_{i=1}^{n} x_i^2}$$

$$= \frac{\left[ \sum_{i=1}^{n}(X_i - \bar{X})^2 - n\bar{X}^2 \right]}{n\sum_{i=1}^{n} x_i^2} \quad \frac{\left[ \sum_{i=1}^{n} X_i^2 - 2\bar{X}\sum_{i=1}^{n} X_i + n\bar{X}^2 - n\bar{X}^2 \right]}{n\sum_{i=1}^{n} x_i^2}$$

$$\frac{\sum_{i=1}^{n} X_i^2 - 2n\bar{X}^2 - 2n\bar{X}^2}{n\sum_{i=1}^{n} x_i^2} = \frac{\sum_{i=1}^{n} X_i^2}{n\sum_{i=1}^{n} x_i^2}$$

Here $\sigma_u$ represents the variance of the unobservable disturbances which is not known. In particular, $\sigma_u$ is not known and we substitute by its unbiased estimator

$$\sum_{i=1}^{n} \ldots = E(\sigma_u) = \sum \ldots (n\ldots)^{-1} = \sigma_u, \quad \text{then the standard normal variable } z$$

or $z$ will follow a $t$-distribution with $(n-2)$ degrees of freedom.

$$\text{In case of } \alpha \quad z = \frac{\alpha - \alpha}{\sigma_u \sqrt{\sum_{i=1}^{n} x_i^2}} = \frac{\sqrt{n \sum x_i}}{\sigma_u \sqrt{\frac{\sum}{n}} \, 1}$$

$$\sigma_u \sqrt{n \sum_{i=1}^{n} x_i^2}$$

When $\sigma_u$ is not known and it is replaced by $\hat{\sigma}_u$, $\hat{\sigma}_u = \sqrt{\left[\sum_{i=1}^{n} e_i^2 \right/ (n-2)\right]}$ the unbiased estimator of $\sigma_u$, then we have.

$$t = t_{n-2} = \frac{\alpha - \alpha}{\hat{\sigma}_u \sqrt{\sum_{i=1}^{n} x_i^2}} \sqrt{n \sum_{i=1}^{n} x_i^2} \quad \text{with d.f} = (n-2)$$

Now by rearranging in terms of $t$ expression we have

$$\alpha - \alpha = t_{n-2} \hat{\sigma}_u \sqrt{\frac{\sum x_i^2}{n \sum_{i=1}^{n} x_i^2}} \quad \text{or} \quad \alpha = \alpha \pm t_{n-2} \hat{\sigma}_u \sqrt{\frac{\sum x_i^2}{n \sum_{i=1}^{n} x_i^2}}$$

Therefore 95% confidence limits for $\alpha$ are

$$\alpha \pm t_{0.025, \, n-2} \, \hat{\sigma}_u \sqrt{\sum_{i=1}^{n} x_i^2 \left/ n \sum_{i=1}^{n} x_i^2 \right.}$$

Similarly, 99% confidence limits for $\alpha$ are

$$\alpha \pm t_{0.005, \, n-2} \, \hat{\sigma}_u \sqrt{\sum_{i=1}^{n} x_i^2 \left/ n \sum_{i=1}^{n} x_i^2 \right.}$$

[The values of $t_{0.025, \, n-2}$ and $t_{0.005, \, n-2}$ corresponding to $(n-2)$ d.f can be obtained from the table, given at the end of the book.]

In the same way, for testing $\beta$ we have

$$t = \frac{\hat{\beta} - \beta}{\sigma_u \sqrt{\dfrac{1}{\sum\limits^{n} x_i^2}}}$$

when $\sigma_u$ is not known then it is replaced by its unbiased estimator $\hat{\sigma}_u$ then we have

$$t = t_{n-2} = \frac{\hat{\beta} - \beta \sqrt{\sum\limits^{n} x_i^2}}{\hat{\sigma}_u} \quad \text{with d.f} = n - 2. \text{ Now rearranging we may get}$$

$$\hat{\beta} - \beta = t_{n-2} \frac{\hat{\sigma}_u}{\sqrt{\sum\limits_{i=1}^{n} x_i^2}} \quad \text{or} \quad \beta = \hat{\beta} \pm t_{n-2} \frac{\hat{\sigma}_u}{\sqrt{\sum\limits^{n}_{1} x_i^2}}$$

Therefore 95% confidence limits for $\beta$ would be

$$\hat{\beta} \pm t_{0.025, n-2} \frac{\hat{\sigma}_u}{\sqrt{\sum\limits_{i=1}^{n} x_i^2}}$$

and 99% confidence limits for $\beta$ would be

$$\hat{\beta} \pm t_{0.005, n-2} \frac{\hat{\sigma}_u}{\sqrt{\sum\limits_{i=1}^{n} x_i^2}}$$

**Confidence interval for $\sigma_u^2$ :**

Under the normality assumption, the variable $\chi^2 = (n-2)\dfrac{\hat{\sigma}_u^2}{\sigma_u^2}$ follows a $\chi^2$ distribution with d.f $= (n-2)$.

Therefore, we can use $\chi^2_{n-2}$ to establish a confidence interval for $\sigma_u^2$

$$P\left[ \chi^2_{1-\frac{\alpha}{2}} \le \chi^2_{n-2} \le \chi^2_{\frac{\alpha}{2}} \right] = 1 - \alpha$$

or, $$P\left[ \chi^2_{1-\frac{\alpha}{2}} \le (n-2)\frac{\hat{\sigma}_u^2}{\sigma_u^2} \le \chi^2_{\frac{\alpha}{2}} \right] = 1 - \alpha$$

Since the statistical test procedure ... the null hypothesis that ... the parameter $\beta$ ...

... test our null hypothesis $H_0: \beta = \beta_0$ against the alternative hypothesis $H_1: \beta \neq \beta_0$ ... Hence the appropriate test statistic under $H_0$ ...

$$t = \frac{\hat{\beta}}{SE(\hat{\beta})} = \frac{\hat{\beta}}{\dfrac{\sigma_u}{\sqrt{\sum x_i^2}}}$$

which follows a $t$-distribution with d.f. $= (n-2)$.

At the level of significance the null hypothesis will be rejected for the given sample if $t_{n-2}$ (observed) $> t_{0.025, n-2}$ and will be accepted otherwise ... Similarly at 1% level of significance the null hypothesis will be rejected for the given sample if $t_{n-2}$ (observed) $> t_{0.005, n-2}$ and will be accepted otherwise (i.e. $-t_{0.005, n-2} < t < t_{0.005, n-2}$). The confidence limits for $\beta$ (acceptance region) in a two tailed test at 5% and 1% levels of significance with $(n-2)$ degrees of freedom will be given by,

$$-t_{0.025, n-2}\, SE(\hat{\beta}) < \beta < +t_{0.025, n-2}\, SE(\hat{\beta})$$

and

$$-t_{0.005, n-2}\, SE(\hat{\beta}) < \beta < +t_{0.005, n-2}\, SE(\hat{\beta})$$

where $SE(\hat{\beta}) = \dfrac{\sigma_u}{\sqrt{\sum\limits_{i=1}^{n} x_i^2}}$ ($\sigma_u$ is not known and replaced by $\hat{\sigma}_u$)

### 2.13.1. The Exact Level of Significance – The p-value

We know that the significance level $\alpha$ in a hypothesis testing problem is the probability of making a Type I error, i.e. $\alpha$ is the probability of rejecting the null

hypothesis ... when ... ... ... ... ... ... ... ... ... ... the maximum ...
... ... the probab... ... ... ... ... ... ... to mak...
... ... of the appropriate ... ... ... ... ... ... ... ...
... ... ... ... ... ... ... ... ... ... ... the testing ... ...
hypoth... the ... ... ... ... the closer ... ... ... ... ... ... ... the
... ... ... ... the tolerance ... ... ... ... ... ... ... ... ... ... then
... rejected and ... otherwise accepted. So also the ... ... the more ... the ...
limit for the observed test statistic gets.

Now we can ask a natural question: what is the smallest value of ... significance
level at which the null hypothesis gets rejected? The answer is p-value... probability
value associated with the observed data set.

Once the p-value is computed we can make use of ... value for ... test and
comparing the p-value with $\alpha$.

**Usually if the p-value $\leq \alpha$ then reject $H_0$, otherwise accept $H_0$.**

To illustrate it we consider an example where we have a linear regression equation of
... ... showing the impact of education on wages given a sample ...
$n$ ... where $Y$ = wages, $X$ = Education years (showing the original data is not
given here).

The estimated regression results are given below

$$Y = \alpha + \beta X$$

$$\Rightarrow \quad Y = -10.44 + 0.7240 X$$

$$SE = (197.2) \quad (0.0700), \quad r^2 = 0.9065$$

Suppose we like to test the null hypothesis $H_0: \beta = 0.5$ against the alternative
hypothesis $H_1: \beta \neq 0.5$. The appropriate test statistic under $H_0: \beta = 0.5$ will be

$$t_{N-2} = \frac{\hat{\beta} - \beta}{SE(\hat{\beta})} = \frac{0.7240 - 0.5}{0.0700} = 3.2$$

$$t_{N-2} \text{ (observed)} = 3.2$$

Now on the basis of the given sample $H_0: \beta = 0.5$ will be rejected at ... level
of significance when $(H_1: \beta \neq 0.5)$ if $t_{\alpha}$ (observed) $> t_{\alpha, n}$ Table

Here $n = 13$, if $\alpha = 0.05$, $t_{\alpha/2, n-2} = t_{0.025, 11} = 2.201$

and if $\alpha = 0.01$, $t_{\alpha/2, n-2} = t_{0.005, 11} = 3.106$

So, $H_0: \beta = 0.5$ is rejected both at 5% and 1% levels of significance as $t_N$
(observed) $= 3.2 > 2.201$ and $3.106$.

Now on the basis of p-value the null hypothesis will be rejected for the given sample
if $p \leq \alpha$ and will be accepted otherwise.

Here given the null hypothesis, that the true coefficient of education $\beta = 0.5$, we
obtain a t value of 3.2. Now what is the p-value of obtaining a t value of as much as
or greater than 3.2 ?

From the t table given in Appendix Table IX we see that for ... d.f. the probability
of obtaining such t value must be smaller than 0.005 (one tail) or 0.01 (two ...

*[Several lines of text at the top of the page are too faded and degraded to read reliably.]*

The case ... of significance ... the $t$ statistic is much smaller than ... an conventional and arbitrary fixed level of significance such as ... if it per ...

As a matter of fact ... we were to use the $p$-value just computed and ... the null hypothesis that the true coefficient of education is ... the probability of committing a type I error would be only about 1 in ...

For a given sample size as ... increases the $p$-value decreases and we can therefore reject the null hypothesis with increasing confidence.

What is the relationship of the $p$-value to the level of significance ... If we make the habit of fixing ... equal to $p$-value of a test statistic (e.g. the ... statistic ... ) there ... distinction between the two values. To put it differently it is the practice ... to ... arbitrary ... at some level and simply choose the $p$-value of the test statistic ... so that we reject the null hypothesis, $H_0$, on the basis of the given sample if $p$-value ... and accept $H_0$ otherwise.

## 2.14 Goodness of Fit of the Multiple Correlation Coefficient ($R^2$)

So far we were concerned with the estimation and precision of the regression parameters $\alpha$ and $\beta$. We now like to consider the regression line as a whole and examine its goodness of fit. Suppose a sample regression has been obtained by the method of least squares, as shown in the following diagram (Fig. 2.5). Considering a specific observation of the dependent variable $Y_i$, we can write $e_i = Y - \hat{Y}$ where $Y = \alpha + \beta X + e_i$ and $\hat{Y} = \alpha + \beta X$, $e_i$ being the error of estimate.



Fig. 2.5

Now $Y_i = \hat{Y} + e_i$

or, $Y_i - \bar{Y} = (\hat{Y}_i - \bar{Y}) + e_i$     $\bar{e} = 0$ or, $(Y_i - \bar{Y}) = (\hat{Y}_i - \bar{Y}) + e_i$     ( $\bar{e} = 0$ )

or $\sum_{i=1}^{n}(Y_i - \bar{Y}) = \sum_{i=1}^{n}(\hat{Y}_i - \bar{Y}) + \sum_{i=1}^{n} e_i$ and $\sum_{i=1}^{n}(Y_i - \bar{Y})^2 = \sum_{i=1}^{n}(\hat{Y}_i - \bar{Y})^2$

$$\text{or} \quad \sum_{i=1}^{n} y_i^2 = \beta^2 \sum_{i=1}^{n} x_i^2 + \sum_{i=1}^{n} e_i^2 \quad \text{where } y_i = Y_i - \bar{Y} \text{ and } x_i = X_i - \bar{X}$$

$$\Rightarrow \begin{bmatrix} \text{Total sum of} \\ \text{squares} \end{bmatrix} = \begin{bmatrix} \text{Explained sum} \\ \text{of squares} \end{bmatrix} + \begin{bmatrix} \text{Unexplained} \\ \text{sum of squares} \end{bmatrix} \text{ or } \begin{bmatrix} \text{Residual sum} \\ \text{of squares} \end{bmatrix}$$

$$\Rightarrow \text{TSS} = \text{ESS} + \text{RSS}$$

$\sum_{i=1}^{n}(Y_i - \bar{Y})^2$ represents the total sum of squared deviations from $\bar{Y}$ which we may take as a measure of the total variations in $Y$.

Thus total variations can be decomposed into two parts

(i) $\beta^2\left[\sum_{i=1}^{n}(X - \bar{X})^2\right] = \hat{\beta}^2 \sum_{i=1}^{n} x_i^2 \Rightarrow$ the estimated effect of change in $X$ on the variations in $Y$ (Explained sum of squares).

Since $\text{var } Y = \text{var } \hat{Y} + \text{var } e$

and $0 \le \text{var } \hat{Y} \le \text{var } Y$

or $0 \le \dfrac{\text{var } \hat{Y}}{\text{var } Y} \le 1$ or $0 \le R^2 \le 1$

$R^2 = 0$ when $\text{var } \hat{Y} = 0$ i.e. $\sum_{i} \hat{e}_i^2 = \sum_{i=1} y_i^2$

and $R^2 = 1$ when $\text{var } \hat{Y} = \text{var}(Y)$ i.e. $\sum_{i} \hat{e}_i^2 = 0$

It should be noted that, this $R^2$ is equal to the square of the simple correlation coefficient between $X$ and $Y$.

By definition simple correlation coefficient (product moment) given by

$$r_{xy} = r = \dfrac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \dfrac{\frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X})^2}\sqrt{\frac{1}{n}\sum(Y_i - \bar{Y})^2}}$$

$$r = \dfrac{\sum_{i=1}^{n} x_i y_i}{\sqrt{\sum_{i=1}^{n} x_i^2 \, \sum_{i=1}^{n} y_i^2}}$$

Since $\beta = \dfrac{\sum_{i=1}^{n} x_i y_i}{\sum_{i=1}^{n} x_i^2}$ and $R^2 = \dfrac{\beta \sum_{i=1}^{n} x_i^2}{\sum_{i=1}^{n} y_i^2}$

$$R^2 = \left[\frac{\sum_{i=1}^{n}}{\sum} \bigg| \frac{\sum_{i=1}^{n}}{\sum} \cdot \frac{\sum}{\sum} = \frac{\sum}{\sqrt{\sum} \cdot \sqrt{\sum}}\right] r$$

$R = r^2$

Since $0 \le R^2 \le 1$

$0 \le r^2 \le 1$

or $-1 \le r \le +1$

**Example 2.4.** Find the value of $R^2$ from the following information and coefficient

$$\sum_{i=1}^{n} = 3347.60, \quad \sum_{i=1}^{n} x_i^2 = 604.80, \quad \sum_{i=1}^{n} y_i^2 = 19837, \quad n = 20, \text{ where } \quad X \quad Y \text{ and }$$

$\bar{Y}$

**Solution** Since $R^2 = \dfrac{\hat{\beta}^2 \sum_{i=1}^{n} x_i^2}{\sum y^2}$ where $\beta = \dfrac{\sum_{i=1}^{n} x_i y_i}{\sum x_i^2}$

$$\beta^2 = \frac{\left(\sum_{i=1}^{n} x_i y_i\right)^2}{\sum x_i^2} = \left(\frac{3347.60}{604.80}\right)^2 \quad (5.54)^2 \quad 30.69$$

Now $R^2 = \dfrac{\hat{\beta}^2 \sum_{i=1}^{n} x_i^2}{\sum y_i^2} = \dfrac{30.69 \cdot 604.80}{19837} = \dfrac{18561.3\,2}{19837} = 0.935$

$R^2 = 0.935$

This suggests that 93.5 percent of the changes in the sample observations of $Y$ can be attributed to the variations of the fitted value of $Y$ i.e. $\hat{Y}$ or we say that our regression line fits the given data well.

Thus $R^2$ measures the proportion of variations in the dependent variable that is explained by the independent variables.

**Example 2.5.** A sample of 20 observations corresponding to the regression model $Y_i = \alpha + \beta X_i + u$ where $u$ is normally distributed with mean zero and unknown variance $\sigma_u^2$ gives the following data .

$$\sum = 7,9 \quad \sum_{} Y \quad \bar{Y}) = 80 \vee \quad \sum (1 \quad \bar{Y} 0), \quad \hat{Y}) = 106.4$$

$$\sum Y = 190.2 \quad \sum (1 \quad \bar{Y} = 7.54 \quad n = 20$$

Obtain the usual regression results.

**Solution** On the basis of the given information we have to fit a linear relation between $Y$ (dependent variable) and $X$ (explanatory variable)

i) Estimation of $\alpha$ and $\beta$

We know that, $\beta = \dfrac{\sum x_i y_i}{\sum x_i^2} = \dfrac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2}$

$$\hat{\beta} = \dfrac{106.4}{2 \, 54} = 0.494$$

and $\alpha = \bar{Y} - \beta \bar{X}$ where $\bar{Y} = \dfrac{\sum Y}{n} = \dfrac{21.9}{20} = 1.095$ and $\bar{X} = \sum \dfrac{X}{n} = \dfrac{186.2}{20} = 9.31$

$= 1.95 \quad 0.494 \times 9.31$

$= 1.95 \quad 4.60 \quad 3.505$

Thus we have, $\alpha = 3.505$ and $\beta = 0.494$ and our estimated regression line is

$\hat{Y} = \alpha - \beta X \Rightarrow Y = 3.505 + 0.494 X$

ii) Estimation of variances

Since we know that, $\text{var}(\alpha) = \sigma_u^2 \left\{ \dfrac{\sum x_i^2}{n \sum x_i^2} \right\}$ and $\text{var}(\beta) = \dfrac{\sigma_u^2}{\sum x^2}$

Here we see that $\sigma_u^2$ is not known and hence we replace it by its unbiased estimator

$$\hat{\sigma}_u^2 = \sum_{i=1}^{n} e_i^2 \quad (n-2)$$

Thus we have, $\text{var}(\alpha) = \hat{\sigma}_u^2 \dfrac{\sum x_i^2}{\left( n \sum x_i^2 \right)}$ and $\text{var}(\beta) = \dfrac{\sigma_u^2}{\sum x^2}$

Again we know that $\sum$ , $\sum$ , $\sum$ ,

Since $\sum$ , $\sum$ , $\sum$ , where $\sum$ , $\sum$ , , , and $\sum$ , $\sum$ , , ,

$\sum$ , 86.9 , (0.494) , 215.4 , 86.9 , 42.56 , 44.34

Now $\sigma_u^2$ , $\dfrac{\sum_{i=1}^{n} e^2}{(n-2)} = \dfrac{44.34}{20-2} = \dfrac{44.34}{n} = 1.908$

Now $\text{var}(\alpha) = \sigma_u$ $\dfrac{\sum x_i^2}{\left(n \sum_i^n e_i^2\right)}$ $= \dfrac{1.908 \cdot 1948.922}{20 \cdot 215.4}$ , 0.8631

$\left[ \sum_{i=1}^{n}(X-X)^2 = 213.4 \text{ or } \sum_{i=1}^{n} x_i^2 - n\bar{x}^2 = 215.4 \right.$

or, $\sum_{i=1}^{n} X_i^2 = 2.54 + n\bar{X}^2 = 215.4 + 20 \cdot 9.3^2$

$= 2.54 + 1733.522 = 1948.922]$

$\therefore \text{var}(\hat{\alpha}) = 0.8631$

Similarly, $\text{var}(\hat{\beta}) = \dfrac{\sigma_u^2}{\sum_{i=1}^{n} x_i^2} = \dfrac{1.908}{215.4} = 0.0089$

Now, $SE(\alpha) = \sqrt{\text{var}(\alpha)} = \sqrt{0.8631} = 0.929$

$SE(\hat{\beta}) = \sqrt{\text{var}(\beta)} = \sqrt{0.0089} = 0.094$

**( ) Construction of confidence intervals**

Now we like to set up a confidence interval for $\alpha$ and $\beta$ at (a) $P = 0.95$ i.e. 5% level of significance) and (b) $P = 0.99$ (i.e. 1% level of significance)

In other words, we like to find the value of 't' that cuts off (a) 0.025 and (b) 0.005 of the area at the tail ends of the distribution on both sides. From table value $t_{0.025}$, $(n-2) = t_{0.025}$, $18-2$ 01 and $t_{0.005}$, $(n-2) = t_{0.005}$, $18$ · 2.878

Therefore 95% confidence interval for $\alpha$ are $\hat{\alpha} \pm t_{0.025}$ $(n-2)$ $SE(\hat{\alpha})$ i.e., $P[\hat{\alpha}$ $t_{0.025}$ $(n-2)$ $SE(\alpha) \leq \alpha \leq \alpha + t_{0.025}$ $(n-2)$ $SE(\alpha)] = 0.95$ and 99% confidence interval for $\alpha$ would be $\hat{\alpha} + t_{0.025}$ $(n-2)$ $SE(\alpha)$.

...

Thus the 95% confidence interval for $\beta$ would be ...

where $\hat{\beta} = 0.494$

$\text{SE}(\hat{\beta}) = 0.094$

**Hypothesis testing** Suppose we like to test the null hypothesis $H_0: \beta = 0$ against the alternative hypothesis $H_1: \beta \neq 0$. Now on the basis of the $p$ ... $H_0: \beta = 0$ be rejected at 5% level of significance if

$$t_n = \frac{\hat{\beta}}{\text{SE}(\hat{\beta})} > \text{observed} > t_{0.025}\ (n-2)\ \text{table value}$$

and ... be accepted otherwise.

Here $t_n = \dfrac{\hat{\beta}}{\text{SE}(\hat{\beta})} = \dfrac{0.494}{0.094} = 5.255$ (where $n = 20$)

Thus we see that $t_{cal} = 5.255 > t_{0.025}\ 18 (=2.101)$ and hence $H_0: \beta = 0$ is rejected (alternative $H_1: \beta \neq 0$ is accepted) at 5% level of significance. So the hypothesis of no relationship between $X$ and $Y$ is to be rejected at 5% level of significance. Similarly it can be tested for 1% level of significance.

## 2.5 Results of Regression Analysis

The results of regression analysis are generally presented in a conventional ... It is not sufficient merely to report the estimates of $\alpha$ and $\beta$. In practice we ... regression coefficients together with their standard errors and the value of $R^2$ ... become customary to present the estimated equation with standard errors put in parentheses below the estimated parameter values. These results are supplemented by $R^2$, the value of which is written on right hand side of the estimated regression equation.

In terms of our earlier example (Example 2.5) the estimated regression results can be written as

$$Y = 7.505 + 0.494\ X \qquad R^2 = 0.6048$$
$$\text{SE} \quad (0.929)\ (0.094)$$

Here ...

$$r \frac{\sum ...}{...}$$

$$... \frac{...}{\sum ...}$$

This suggests that variations in ... percent ... be attributed to the variations of the total ... it is ... for the given data moderate ... not very ...

Some economists and report the ... of the estimated ... array of standard errors. This way of presentation makes the testing of hypothesis easier and direct.

Thus the other form of presentation of results in the ... example ...

$$... \quad ... \quad ... \quad R^2 = ...$$

... (... ), (... )

where $\alpha = 1.505 \quad \beta = 0.494$

$$\frac{\alpha}{SE(\alpha)} = ... \quad \frac{\beta}{SE(\beta)} = 5.255$$

**Example 2.6.** Suppose that Mr. X estimates a consumption function and shares the results

$$C = \quad 5 + 0.8 \, Y_d \qquad n = 19$$
t ratios $\quad (3.1) \ (18.7) \qquad R^2 = 0.99$

$C$ is consumption, $Y_d$ is disposable income the numbers in parentheses are t ratios

(a) Test the significance of $Y_d$ statistically using t-ratios

(b) Determine the estimated standard derivations of the parameter estimators

(c) Construct a 95 percent confidence interval for the coefficient of $Y_d$

**Solution :** It is a formal consumption function of the Keynesian type $C = a + bY_d$ where $a$ = autonomous part of consumption and $b$ = Marginal propensity to consume. By assumption in the existing theory $a > 0$, $b < 1$. The estimated relation/regression results are given here as

$$\hat{C} = \quad 15 + 0.81 \, Y_d \qquad n = 19$$
t ratios $\quad (3.1) \ (18.7) \qquad R^2 = 0.99$

This shows that $\hat{a} = 15$, $\hat{b} = 0.81$,

$$\frac{\hat{a}}{SE(a)} = 3.1 \ (t\text{-ratio}) \quad \text{and} \quad \frac{\hat{b}}{SE(b)} = 18.7 \ (t\text{-ratio})$$

$n$ = no. of data points (sample size = 19)

$R^2$ = Square of multiple correlation coefficient

$$= \frac{ESS}{TSS} = \frac{\text{Explained variation in } C}{\text{Total variation}} = \frac{\text{var}(\hat{C})}{\text{var}(C)}$$

where $b$ = ...

where $\sum \quad \sum \quad + \sum$

Here $r^2$ = ...

This means that 99 percent of the variations in sample observations of ... is attributed to the variations of the fitted value of ... Thus we say that ... regression line on the given data well. That of 100%, variation in consumption ...

... regression relation can explain 99% variation in consumption.

(a) We have to test the null hypothesis $H_0$ : $b = 0$ no relation between $C$ and ... against the alternative hypothesis $H_1$ : $b \neq 0$.

The appropriate test statistic under $H_0$ : $b = 0$ would be

$$\frac{b}{SE(b)} \quad \text{which follows a } t\text{-distribution with } (n-2) \text{ degrees of freedom}$$

i.e. $t_b = t = \frac{b}{SE(b)} = 18.7 \text{ (given)}$

Now at 5% level of significance $H_0$ : $b = 0$ (no relation between $C$ and $Y_d$)

will be accepted if $-t_{0.025} \leq t \leq t_{0.025} \quad df = 2$

and will be rejected otherwise.

From table value we get

$t_{0.025} = t_{0.025, 17} = 2.110$

$(n = 19 \text{ given})$

Thus we see that observed $t = \frac{b}{SE(b)} = 18.7$ does not lie in the interval $-2.110$ and $2.110$ and hence the null hypothesis is rejected and the alternative is accepted. This means that there exists a relation between consumption (C) and disposable income $Y_d$. Hence the relation is statistically significant.

(b) We have to find $SE(a)$ and $SE(b)$

Since for $a, t = \frac{a}{SE(a)} = 3.1 \text{ (given)}$

and $a = 15$, $3.1 = \frac{15}{SE(a)}$ or, $SE(a) = \frac{15}{3.1} = 4.8387$

Similarly, for $b$, $t = \frac{b}{SE(b)} = 18.7 \text{ (given)}$

and $b = 0.81$

or, $18.7 = \frac{0.81}{SE(b)}$ or, $SE(b) = \frac{0.81}{18.7} = 0.0433$

Thus the estimated standard deviations of the parameter estimators are

...

(i.e. β)

This will be given by

...

e. between ( ... ) and ( ... )

This means that the coefficient of $Y$,
will be between 0.7187 and 0.9013

**Example 2.7** The following table shows data on Labour hours of work and output for ... workers.

| X (labour hours of work) | 10 | 7 | 10 | 5 | 8 | 8 | 6 | 7 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Y (output) | 11 | 10 | 12 | 6 | 10 | 7 | 9 | 10 | 11 | 10 |

Assuming a linear regression of the form $Y = \alpha + \beta X$ ...

the OLS estimators of $\alpha$ and $\beta$ (i.e. $\hat\alpha$ and $\hat\beta$) and the estimated regression line $\hat Y = \hat\alpha + \hat\beta X$

(ii) Find var($\hat\alpha$), var($\hat\beta$), SE($\hat\alpha$) and SE($\hat\beta$)

(iii) Find the value of $\sum_{i=1}^{n} e_i^2$ and $\hat\sigma_u^2 = \dfrac{\sum_{i=1}^{n} e_i^2}{n-2}$

(iv) Find the value of $R^2$

(v) Construct 95% confidence intervals of $\alpha$, $\beta$ and $\sigma_u^2$

(vi) Test the null hypothesis $H_0: \beta = 1.35$ against
a) $H_1: \beta = 1.35$, (b) $H_1: \beta > 1.35$ (c) $H_1: \beta < 1.35$

**Solution :**      **Calculations for the Regression**

| Observations | $X_i$ | $Y_i$ | $x_i = X_i - \bar X$ | $y_i = Y_i - \bar Y$ | $\hat Y_i$ | $e_i$ | $x_i y_i$ | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 10 | 11 | 2 | 1.4 | 4 | 2.8 | | | | |
| 2 | 7 | 10 | 1 | 0.4 | 1 | 0.4 | 8.45 | | | 5 |
| 3 | 10 | 12 | 2 | 2.4 | 4 | 4.8 | 11.5 | | | 6.9 |
| 4 | 5 | 6 | 3 | 3.6 | 9 | 10.8 | 7.45 | | | 4.5 |
| 5 | 8 | 10 | 0 | 0.4 | 0 | 0 | 9.60 | | | 4 |
| 6 | 8 | 7 | 0 | 2.6 | 0 | 0 | 9.60 | | | 3.6 |
| 7 | 6 | 9 | 2 | 0.6 | 4 | 1.2 | 8.0 | | | 1.9 |
| 8 | 7 | 10 | 1 | 0.4 | 1 | 0.4 | 8.65 | | | 1.5 |
| 9 | 9 | 11 | 1 | 1.4 | 1 | 1.4 | 10.35 | | | 6.5 |
| 10 | 10 | 10 | 2 | 0.4 | 4 | 0.8 | 0 | | | |

| Total | $\sum_{i=1}^{n} X_i$ | $\sum_{i=1}^{n} Y_i$ | $\sum_{i=1}^{n} x_i$ | $\sum_{i=1}^{n} y_i$ | $\sum_{i=1}^{n} x_i^2$ | $\sum_{i=1}^{n} x_i y_i$ | $\sum_{i=1}^{n} \hat Y_i$ | $\sum e_i = 0$ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | = 80 | = 96 | = 0 | = 0 | = 28 | = 21 | = 96 | | | |

Here n ... $\bar{x} = \sum ...$ ... $n = ...$ ... $\bar{y} = \sum ...$ ... $\frac{y...}{n} ...6$

$$\sum ... ... ... ... = 11.65$$

and $\sum ...$ ... ... ... = ...668

... the OLS estimators of $\alpha$ and $\beta$ are given by $\hat{\alpha}$ and $\hat{\beta}$ where

$$\beta = \sum ... \bigg/ \sum ... = \frac{...}{...} = 0.75$$

and ... $\bar{y} - \beta \bar{x} = 9.6 - 0.75 ... = 9.6 - ...6$

$$\alpha = 3.6 \text{ and } \beta = 0.75$$

The estimated regression line equation becomes

$\hat{y} = \alpha + \beta x$ or $\hat{y} = 3.6 + 0.75 x$. This equation is now used to find the values $\hat{y}$ corresponding to different values of $x$. The values of $\hat{y}$ are given in the above table showing calculations for regression.

Here the regression coefficient $\beta = 0.75$ measures the marginal product ivity of labour. The intercept $\alpha = 3.6$ means that output will be 3.6 units when ... hours of work is zero.

(ii) We have to find $\text{var}(\alpha)$, $\text{var}(\beta)$, $SE(\alpha)$ and $SE(\beta)$.

Since we know that $\text{var}(\alpha) = \sigma_u^2 \cdot \dfrac{\sum x^2}{n \sum x^2}$ and $\text{var}(\beta) = \sigma_u^2 \cdot \dfrac{1}{\sum x^2}$

Since $\sigma_u^2$ is unknown it is replaced by its unbiased estimator

$$\sigma_u^2 = \sum_{i=1}^{n} \hat{u}_i^2 /(n-2) = \frac{14.65}{10-2} = \frac{4.65}{8} = 1.8312$$

Now $\text{var}(\alpha) = \dfrac{\sigma_u^2 \sum_{i=1}^{n} x}{n \sum_{i=1}^{n} x} = \dfrac{... \times 668}{10 \times ...} = \dfrac{...416}{...} = 4.5687$

and $SE(\alpha) = \sqrt{\text{var}(\alpha)} = \sqrt{4.5687} = ...$

Again ...

We ... of the ... $\sum$ ...

... $\sum$ ...

(taken from last column of calculation table)

and ... $\sum$ ...

(iv) We have to find out the value of the ... R

Since we know that $R^2 = \dfrac{ESS}{TSS} = \dfrac{\hat{\beta}^2 \sum\limits_{i=1}^{n} x_i^2}{\sum\limits^{n} y_i^2}$    (Since $\sum\limits_{i=1}^{n} y_i^2 = \hat{\beta}^2 \sum\limits_{i=1}^{n} x_i^2 + \sum\limits$

i.e., $TSS = ESS + RSS$)

$$\sum\limits_{i=1}^{n} y_i^2 = \hat{\beta}^2 \sum\limits_{i=1}^{n} x_i^2 + \sum\limits_{i=1}^{n} e_i^2$$

$$= (0.75)^2 \times 28 + 14.65 = 15.75 + 14.65 = 30.40$$

$$R^2 = \dfrac{\hat{\beta}^2 \sum\limits_{i=1}^{n} x_i^2}{\sum\limits^{n} y_i^2} = \dfrac{(0.75)^2 \times 28}{30.40} = \dfrac{15.75}{30.40} = 0.51$$

This suggests that 51 percent of the variations in the sample observations (Y) can be attributed to the variations of the fitted value of Y i.e. $\hat{Y}$. Here we see that our regression line fits the given data moderately (not very well).

From the above results we can write our regression results as follows:

$$\hat{Y} = 3.6 + 0.75\, X,\quad R^2 = 0.51$$

(2.090)  (0.256)   [SE values in brackets]

Alternatively, $\hat{Y} = 3.6 + 0.75\, X,\quad R^2 = 0.51$

(1.7224) (2.930)   [t values in brackets]

(v) 95% confidence intervals of $\alpha$, $\beta$ and $\sigma_u^2$

(a) 95% confidence interval for $\alpha$ is given by,

$P[\hat{\alpha} - t_{0.025, n-2}\, SE(\hat{\alpha}) \le \alpha \le \hat{\alpha} + t_{0.025, n-2}\, SE(\hat{\alpha})] = 0.95$

95% confidence limits of $\alpha$ are

Since

from table value,

So 95% confidence limits for $\alpha$ are 1.77 and 8.42

(b) 95% confidence interval for $\beta$ is given by

$$\hat\beta - SE(\hat\beta) < \beta < \hat\beta + SE(\hat\beta)$$

95% confidence limits of $\beta$ are

$$\hat\beta \pm t_{0.025} \cdot SE(\hat\beta)$$

or $0.75 \pm 2 \times 0.256$

Since $t_{0.025,2} = t_{0.025,8} = 2.06$

(From table value)

$= 0.6$ and $1.34$

95% confidence limits of $\beta$ are 0.16 and 1.34

(c) Since we know that $1 - \alpha$ = 95% confidence interval for $\sigma_u^2$ is given by

$$P\left[ (n-2) \frac{\hat\sigma_u^2}{\chi^2_{\alpha/2}} \le \sigma_u^2 \le (n-2) \frac{\hat\sigma_u^2}{\chi^2_{1-\alpha/2}} \right] = 1 - \alpha$$

Here $n = 10$, $n-2 = 0.05$, $\hat\sigma_u^2 = \sum e_i^2 /(n-2) = 1.8312$, $\chi^2_{0.025,8} = 17.535$ (table value)

and $\chi^2_{0.975,8} = 2.180$ (Table value)

So 95% confidence interval for $\sigma_u^2$ would be

$$P\left[ \frac{14.2}{7.535} \le \sigma_u^2 \le 8 \times \frac{14.2}{2.180} \right] = 0.95$$

or $P[0.84 < \sigma_u^2 < 6.72] = 0.95$

95% confidence limits of $\sigma_u^2$ are 0.84 and 6.72

(iv) To test the null hypothesis $H_0$: $\beta = 1.35$ against the alternative hypothesis $H_1$: $\beta \ne 1.35$ the appropriate test statistic would be, $t = \dfrac{\hat\beta - \beta}{SE(\hat\beta)}$

Now on the basis of the sample data $H_0$: $\beta = 1.35$ will be rejected at 5% level of significance if

$$t_{n-2} = \left| \frac{\hat\beta - \beta}{SE(\hat\beta)} \right| \text{(observed)} > t_{0.025,n-2} \text{ (Table value)}$$

and will be accepted otherwise

here ... $\beta$ $\beta$ 0.75 1.35 0.6

... $-2.343$ and $t_{n-2} = 2.343$

From the table value we ...

This ... ... This means that in the basis of sample data ... $\beta$ ... is rejected at 5% level of significance

(b) The null hypothesis $H_0$: $\beta$ ... will be rejected against the alternative ... at 1000% level of significance if for the given sample

$$t_n \ (\text{observed}) = \frac{\beta - ...}{\sqrt{...}} \qquad t_{a, n} \ \text{table value}$$

and will be accepted otherwise. Here $\alpha = 0.05$, $n = 1)$

$$\text{and } _{n-2} (\text{observed}) = \frac{\beta - \beta}{SE(\beta)} \qquad 2.343 \quad t_{a+1} \quad 1.860)$$

So $H$: $\beta = 1.35$ is accepted (hence in significant) at 5% level of significance.

c The null hypothesis, $H_0 = 1.35$ will be rejected against the alternative $H$: $\beta$ ... for the given sample at 1000% level of significance if

$$t_{n-2} \ (\text{observed}) \quad \frac{\beta - \beta}{SF(\beta)} \qquad t_{a, n-2} \ (\text{table value})$$

and will be accepted otherwise.

$$\text{Here, } t_{n-2} \ (\text{observed}) = \frac{\beta - \beta}{SE(\beta)} \qquad \frac{0.75 - 1.35}{0.256} = -2.343$$

$< t_{0.05, 8} = -1.860$ Here $\alpha = 0.05$, $n = 10$

This clearly shows that the null hypothesis $H_0$: $\beta = 1.35$ is rejected (significant) at 5% level of significance.

## 2.16. Analysis of Variance for the Simple Linear Regression Model

Yet another item that is often presented in connection with the simple linear regression model is the analysis of variance. This is the breakdown of the total sum of squares (TSS) into explained sum of squares (ESS) and the residual sum of squares RSS). The purpose of presenting the table is to test the significance of the explained sum of squares. In this case this amounts to testing the significance of $\beta$.

In regression analysis, we minimise the square deviations from mean and it has been proved (see Section 2.14) that.

$$\sum_{i=1}^{n} (Y_i - \bar{Y})^2 = \sum_{i=1}^{n} (Y_i - \bar{Y}_i)^2 + \sum_{i=1}^{n} (e_i)^2 \quad \text{or} \quad \sum_{i=1}^{n} y_i^2 = \sum_{i=1}^{n} \hat{y}_i + \sum_{i} e_i^2$$

That is, Total variation = Explained variation + Unexplained variation or Residual variance or, TSS = ESS + RSS with degrees of freedom $n - 1 = (K - 1) + (n - K)$ where $n$ = total number of observations (given) and $K$ = number of parameters to be estimated.

Now

$$\sum \qquad \sum \qquad \sum \qquad \sum \qquad + \sum$$

and we have

$$\sum \qquad \sum \qquad \beta \sum \qquad \sum$$

and hence we have

$$\sum \qquad \beta \sum \qquad \sum \qquad \sum = \beta^2 \sum \qquad \sum$$

i.e. TSS ESS RSS

with d.f. $n$ — $k$. Here $k = 2$ as there are two parameters $\alpha, \beta$.

Thus we see that total variations are split into explained by explanatory variable and unexplained effect (error) variations again between and within variation. So we ... analysis of variance procedure. This suggests that we can compute the analysis of variance table for the regression analysis also in order to find the overall significance of the regression results.

**ANOVA TABLE**

| Source of variation | Sum of squares | Degrees of freedom | Mean sum of squares | Observed | Tabulated |
|---|---|---|---|---|---|
| Explained between | ESS $\beta \sum$ | $k$ | ESS/$k$ = MSE | $F = \dfrac{MSE}{MSR}$ with $k$ = $k$ | |
| Residual within | RSS $\sum \hat{e}^2$ | $n$ — $k$ | RSS/$n$ — $k$ MSR | $n$ — $k$ | |
| Total | TSS = $\sum y^2$ | $n$ | | | |

Here $K = 2$ as the model is a two variable regression model and two parameters are involved.

To test the null hypothesis ... that the ... in a ... regression where the explanatory variable is ... ... ... we proceed ... by computing $F$ ratio, through ... sum of squares of explained variation ... that ... in testing $H_0$ ... against the alternative $H_1$ ... we may use the test statistic.

$$F = \frac{MSE}{MSN} = \frac{\hat{\beta}^2 \sum x_i^2 \, (K-1)}{\sum e_i^2 \, (n-K)}$$

$$= \frac{\hat{\beta}^2 \sum x_i^2 \, 1}{\sum e_i^2 \, (n-2)} \quad \text{with } df = 1, (n-2) \text{ since } K = 2$$

Now we have to compare $F^*$ ... with the table value of $F$ with $1, (n-2)$ ... if it is found that $F^* > F_{...}$ (table) we reject the null hypothesis ... at the ... of significance ... $0.01$ or $0.05$) and accept otherwise. The ... will be significant ... $H_0$ is rejected.

It should be noted that

$$F^* = \frac{MSE}{MSR} = \frac{\hat{\beta}^2 \sum x_i^2 \, (n-2)}{\sum e_i^2} = \frac{(n-2)\hat{\beta}^2 \sum x_i^2}{(1-R^2)\sum y_i^2}$$

Now $\dfrac{\hat{\beta}^2 \sum x_i^2}{\sum y_i^2} = \dfrac{\sum x_i y_i}{\sum x_i^2} \cdot \dfrac{\sum x_i y_i}{\sum y_i^2}$

But $\dfrac{\sum e_i^2}{\sum y_i^2} = 1 - R^2$ or, $\sum e_i^2 = (1-R^2)\sum y_i^2$

Therefore, $F^* = \dfrac{(n-2)\left\{\dfrac{\hat{\beta}^2 \sum x_i^2}{\sum y_i^2}\right\}}{1-R^2} = \dfrac{(n-2)R^2}{1-R^2}$

which, on generalisation becomes $\dfrac{R^2 (K-1)}{(1-R^2)(n-K)}$ where we have $K$ parameters

Furthermore, we know that

$$t = \frac{\hat{\beta}}{SE(\beta)} = \frac{\beta}{\sqrt{var(\hat{\beta})}} \quad \text{But } var(\hat{\beta}) = \frac{\sigma_u^2}{\sum x_i^2} = \frac{\sum e_i^2 \,(n-2)}{\sum x_i^2}$$

$$t^2 = \frac{\hat{\beta}^2}{var(\hat{\beta})} = \frac{\beta^2}{\{\sum e_i^2 \,(n-2)\}\left(\dfrac{1}{\sum x_i^2}\right)} \quad \text{or, } t^2 = \frac{\hat{\beta}^2 \sum x_i^2}{\sum e_i^2 \,(n-2)} = F^*$$

The so-called $bt$ and $F$ tests are formally equivalent, the relation between the two being ...

**Example 2.8.** Let us consider the following data to construct the analysis of variance table for a simple regression model $Y_i = \alpha + \beta X_i + u_i$

Given $\sum$ ... $\sum$ ... $\bar{Y}_i = 86.9$ $\sum$ ... ... $Y_i$ ... 186.4

$$\sum i = 80 \sum i \quad 1) = 2 \, 54 \, n = 70$$

**Solution** (See Example 2.3)

The OLS estimators of $\alpha$ and $\beta$ can be obtained as follows

$$\beta = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sum (X_i - \bar{X})^2} = \frac{86.4}{2.54} = 0.494$$

and $\alpha = \bar{Y} - \beta \bar{X}$
$= 0.95 - 0.494 \times 9.31 = 0.95 - 4.60 = 3.505$

where $\bar{Y} = \frac{\sum Y_i}{n} = \frac{21.9}{21} = 0.95$ and $\bar{X} = \frac{\sum x_i}{n} = \frac{186.2}{20} = 9.3$

The estimated regression results are

$\hat{Y} = 3.505 + 0.494 \, X$ $R^2 = 0.6048$

$t$ ratios $(3.772)$ $(5.255)$ where $\alpha = 3.505$, $\beta = 0.494$

$SE(\alpha) = 0.929$ $SE(\beta) = 0.094$ $\dfrac{\alpha}{SE(\alpha)} = 3.772$ $\dfrac{\beta}{SE(\beta)} = 5.255$

$$R^2 = \frac{\beta^2 \sum x_i^2}{\sum y_i^2} = \frac{(0.494)^2 \times 2\,54}{86.9} = 0.6048$$

Now $\sum e_i^2 = \sum y_i^2 - \beta^2 \sum x_i^2 = 86.9 - (0.494)^2 \times 215.4 = 34.14$

Now for testing $H_0$: $\beta = 0$
against the alternative $H_1$: $\beta \neq 0$, we may use the ANOVA Table

## ANOVA TABLE

| Sources of variation | Sum of squares | Degrees of freedom | Mean sum of squares | F ratio | F theoretical at 1% and 5% |
|---|---|---|---|---|---|
| Explained (between) | $TSS$ $\hat{\beta} \sum_{1}^{n}$ $52.56$ | $k$ $= 2$ | $MSE \dfrac{ESS}{k}$ $\dfrac{52.56}{1}$ | $F = \dfrac{MSE}{MSR}$ $\dfrac{52.56}{1.908}$ $= 27.55$ | $F_{0.01}$ $= 8.29$ |
| Residual (within) | $RSS$ $\sum_{1}^{n} e^2$ $34.34$ | $n - k$ $= 20 - 2 = 18$ | $MSR \dfrac{RSS}{n-k}$ $\dfrac{34.34}{18}$ $1.908$ | with df | $F_{0.05\ (1,18)}$ $4.4$ |
| Total | $TSS$ $\sum_{1}^{n} y^2$ $86.9$ | $n - 1$ $= 20 - 1 = 19$ | | | |

Here we see that the observed $F^* = 27.55$ is much larger than table $F_{0.01}$ $= 8.29$ and $F_{0.05}$ $= 1.18 = 4.41$. This means that $H_0$: $\beta = 0$ is rejected both at 5% and 1% levels of significance. Hence we reject the null hypothesis and accept that the regression is significant, that is, $X$ is a significant explanatory factor of the variation in $Y$.

## 2.17  Testing the Equality between Coefficients Obtained from Different Regressions or Different Samples

Sometimes we may have to estimate a regression equation separately for several sets of data and we may have to test whether some or all the parameters are the same for all different sets of data.

Suppose, we have two samples on the variables $Y$ and $X$ containing $n_1$ observations for first set and, $Y$ and $X$ containing $n_2$ observations for second set. We may obtain two estimates of the same relationship for these two samples

$$Y_t = \hat{\alpha}_1 + \hat{\beta}_1 X$$

and   $$Y_t = \hat{\alpha}_2 + \hat{\beta}_2 X$$

Now our problem is to examine whether these two estimated relations differ significantly. If yes, then we may conclude that the relationship changes from one sample to the other.

For example, suppose that we have the data on consumption and disposable income for the two periods 1990-1999 and 2000-2019. We estimate the consumption functions separately for these periods. Then we may be interested to examine whether the functions are statistically significant or whether the MPC significantly differ or not.

Step 1 We have to ... the pooled ... data with number of observations $n = n_1 + n_2$ ... degrees of freedom ... and ... the pooled data.

Step 2 Next ... the regression for each sample separately.

For the sample ... $D_1$ ... and $\sum$ ...

For second sample ... and $\sum c_2 = \sum$ ...

Step 3 Next we have to compute $F$ value as follows

$$F^* = \frac{\sum_1 \quad \sum_2 \quad \sum \quad /k}{\sum - \sum c \quad n_1 + n \quad 2k} \quad \text{with d.f. } k, (n_1 + n_2 - 2k)$$

Next we have to test the null hypothesis

$H_0 : \beta_1 = \beta_2 = \ldots$ against the alternative $H_1$: $H_0$ is not correct

... we reject the null hypothesis at 5% level of significance.

In particular, if the pooled results are not given then $F^*$ value can be obtained as follows

$$F^* = \frac{\sum e^2 \quad /n_1 \quad 2}{\sum e_i^2 \quad n - 2} \quad \text{with d.f. } (n_1 - 2), n_2 \quad ?$$

**Example 2.9** In order to test the null hypothesis that there is no difference in the MPC (Marginal propensity to consume) of manual workers and white-collar employees, a research team estimated the following consumption functions.

**Manual workers** Sample size $n_1 = 35$

$$c_1 = 20 + 1.905 \, y \qquad 0.92 \sum Ic_1 \quad (c_1)^2 = 3.251$$
$$(2) \quad (5.6)$$

(The numbers in brackets are the $t$ values for the regression coefficients)

**White-collar employees** Sample size $n_2 = 30$

$$c_2 = 100 \; 0.825 \quad r^2 = 0.95 \sum Ic \quad \bar{c} \, r = 4.532$$
$$(2) \quad (8.5)$$

The numbers in brackets are the $t$ values for the regression coefficients.

**Combined sample consumption function**
sample size $n = n_1 + n_2 = 30 + 35 = 65$

$$c = 250 + 0.70y \quad r^2 = 0.93 \sum c_p = 6.320$$
$$(5.3) \quad (6.7)$$

On the basis of the above results, can we accept the hypothesis that there is no differences between the MPCs of the two groups? Use a 5 per cent level of significance?

Solution Suppose the estimated consumption function for the first sample is ... and for the second sample ... $+ \beta^{(2)}$

and the pooled estimated consumption function is ... We have to test the ... hypotheses

$H_0$: ... against the alternative $H_1$: $H_1$ is not true

where $\beta^{1}$ = MPC of the first sample

$\beta^{2}$ = MPC of the second sample

$\hat{a}_1$ = MPC of the pooled samples

The appropriate test statistic will be

$$F^* = \frac{\left[\sum c_p^2 - \left(\sum c_1^2 + \sum c_2^2\right)\right] / K}{\left(\sum c_1^2 + \sum c_2^2\right) / (n_1 + n_2 - 2K)} \quad \text{with}$$

$d = [K, (n_1 + n_2 - 2K)]$

Now $H_0$ will be rejected (significant) if $F^* > F_{0.05}, K, (n_1 + n_2 - 2K)$ and will be accepted otherwise.

Now on the basis of our given information we see that

From first sample

$n_1 = 35$, $\eta^2 = R_1^2 = 0.92$, $\sum (c_1 - \bar{c}_1)^2 = \sum c_1^2 = 3,251$

$K = 2$, $\hat{\beta}_1 = 0.90$, $\sum e_1^2 = (1 - R_1^2)\sum c_1^2 = (1 - 0.92) \times 3,251 = 260.08$

$$\left[\text{since } R^2 = 1 - \frac{\sum e^2}{\sum y^2}, \quad \frac{\sum e_1^2}{\sum c_1^2} = 1 - R^2 \text{ or } \sum e_1^2 = (1 - R^2)\sum y^2\right]$$

From second sample $n_2 = 10$, $\beta_2 = 0.82$, $c_2^2 = R_2^2 = 0.95$

$$\sum (c_2 - \bar{c}_2)^2 = \sum c_2^2 = 4,532, \quad K = 2$$

and $\sum e_2^2 = (1 - R_2^2)\sum c_2^2 = (1 - 0.95) \times 4,532 = 226.6$

From the pooled sample $n = n_1 + n_2 = 35 + 10 = 65$

$\hat{a}_1 = 0.70$, $r^2 = R^2 = 0.92$, $\sum c_p^2 = 16,320$

Thus we have, $\dfrac{\left[\sum c_p^2 - \left(\sum c_1^2 + \sum c_2^2\right)\right] / K}{\left(\sum e_1^2 + \sum e_2^2\right) / (n_1 + n_2 - 2K)}$

$$= \frac{[16,320 - (260.08 + 226.6)] / 2}{260.08 + 276.6) / (35 + 10 - 2 \times 2)} \quad \text{with d f} = (K, n_1 + n_2 - 2K)$$

$$= \frac{(16,320 - 486.86) / 2}{486.68 / 61} = \frac{7916.66}{7.9783} = 992.27$$

$F^* = 992.27$ with d f $(2, 61)$.

From the table value we see that $F$ _ _ _ _ 48

Thus we see that _ _ a much larger than the value _ _ of _ $4N$ at _ _ _ _
significance _ hence the null hypothesis is _ he rejected _ _ MPF _ _ the two _ _
_ _ then _ _ _ level of significance We may thus conclude that the MPC _ _
cases differ significantly and _

In particular _ the pooled data are not given then

$$J^* = \frac{\sum_{\cdot} \quad \cdots \quad \cdots}{\sum_{\cdot} \quad \cdots \quad 2} \qquad \text{with d f} \quad \cdots \quad \cdots$$

Here $\cdot^*$ $\dfrac{2 \cdot \quad \cdots \quad 35 \quad 2 \cdot}{2 \quad \cdot \cdot \cdot \quad \cdot)} \quad \text{with d f} (3 + 3) \cdot$

$\dfrac{\cdots \cdots \cdots}{\cdots \cdot 3} \quad \cdots \quad 9.30$

From the table value we see that $F_{0.05}$ _ _ _ $> 1$ 84 (approx

We may thus conclude that the null hypothesis will be accepted $t \approx r$ =
$F_{0.05}$ _ _ at _ % level of significance and hence there will be no difference in MP
in the two cases

## 2.18 Extension of Linear Regression Model to Non-linear Relationships

In the simple linear regression model we consider a linear relation between two
variables $X$ and $Y$ in the form $Y = \alpha + \beta X + u$. But in many situations this may not
be the case In economics we observe non-linear relationships among the variables

Some of the most common forms of non-linear relations used in economics are given
below

Demand curve with unit elasticity $D = f(P)$ or $D = \dfrac{u}{P}$ where $D$ represents
quantity demanded and $P$ denotes price

( ) Average cost curve The traditional theory of ∪ shaped cost curve may be
approximated by a polynomial of third degree in output

$$C = f(q \text{ or } C = \alpha + \beta_1 q_i - \beta_2 q_i^2 + \beta_3 q_i^3 + u_i$$

where $C$ represents cost and $q$ represents the level of output

Now Average cost $\dfrac{C_i}{q_i} = \dfrac{\alpha}{q_i} - \beta + \beta_2 q_i + \beta_3 q_i^2$ which is ∪ shaped curve

( ) Production function may be of the form $Q = f(L, K)$ or $Q = AL^\alpha K^\beta u$ where
$Q$ = level of output $L$ = labour employed $K$ = capital employed $\alpha$ and $\beta$ are two
parameters This type of production function is called Cobb-Douglas production
function

(iv) The production function may be of the form

$$Q_i = A \left[ \delta K_i^{-p} + (1 - \delta) L_i^{-p} \right]^{-\frac{1}{p}}$$

This type of production function is called CES production function The symbols
have their usual meaning

Now to estimate the parameters ... the non linear function ... near form and then we have ... in the usual manner ... method

In order to find out the rate of change ... the slope ... and then we want to find out the elasticity of the regressand with respect to the regressor we use the following table. The knowledge of these ... help us to compose the various models.

| Model | Equation | Slope $\left(\dfrac{dY}{dX}\right)$ | Elasticity $\left(\dfrac{dY}{dX}\cdot\dfrac{X}{Y}\right)$ |
|---|---|---|---|
| Linear | $Y = \alpha + \beta X$ | $\beta$ | $\beta\left(\dfrac{X}{Y}\right)$ * |
| Log linear | $\log Y = \alpha + \beta \log X$ | $\beta\left(\dfrac{Y}{X}\right)$ | $\beta$ |
| Log linear | $\log Y = \alpha + \beta X$ | $\beta(Y)$ | $\beta(X)$ * |
| Linear log | $Y = \alpha + \beta \log X$ | $\beta\left(\dfrac{1}{X}\right)$ | $\beta\left(\dfrac{1}{Y}\right)$ * |
| Reciprocal | $Y = \alpha + \beta\left(\dfrac{1}{X}\right)$ | $-\beta\left(\dfrac{1}{X^2}\right)$ | $-\beta\left(\dfrac{1}{XY}\right)$ * |
| Log reciprocal | $\log Y = \alpha - \beta\left(\dfrac{1}{X}\right)$ | $\beta\left(\dfrac{Y}{X^2}\right)$ | $\beta\left(\dfrac{1}{X}\right)$ * |

**Note** : * indicates that the elasticity is variable, depending on the value taken by $X$ or $Y$ or both. When no $X$ and $Y$ values are specified, in practice, very often these elasticities are measured at the mean values of these variables, namely, $\bar{X}$ and $\bar{Y}$.

**Example 2.10.** Estimate the investment function $I = f(r) = \alpha(r)^\beta u$ on the basis of the following information:

$$n = 11, \quad \sum_{i=1}^{n} Y_i = 12.2771, \quad \sum_{i=1}^{n} X_i = 16.6729$$

$$\sum_{i=1}^{n} X_i^2 = 27.9605, \quad \sum_{i=1}^{n} X_i Y_i = 15.1222,$$

$$\sum_{i=1}^{n}(X_i - \bar{X})(Y_i - \bar{Y}) = 3.4864, \quad \sum_{i=1}^{n}(X_i - \bar{X})^2 = 2.6891,$$

$$\sum_{i=1}^{n}(Y_i - \bar{Y})^2 = 4.8566,$$

where $Y = \log I$, $X = \log r$

**Solution**   The regression function is given by $y = \alpha \beta^x u$ which is a non-linear form where $\alpha$ and $\beta$ are the two parameters whose values are to be estimated by the OLS method. Taking log on both sides we get

$\log y = \log \alpha + x \log \beta + \log u$  or  $y^* = \alpha^* + \beta^* x + u^*$

where $\log y = y^*$, $\log \alpha = \alpha^*$, $\log \beta = \beta^*$ and $\log u = u^*$. This transformation is now linear in terms of logarithms.

The function $y^* = \alpha^* + \beta^* x + u^*$ is of the form $y = \alpha + \beta x + u$ and hence we apply the OLS method.

Now by the OLS method we can obtain the estimates of the parameters $\alpha^*$ and $\beta^*$. Thus we have

(i) $\beta = \dfrac{\sum\limits_{i=1}^{n} X_i - \bar{X})(Y - \bar{Y})}{\sum\limits_{i=1}^{n} X_i - \bar{X})^2} = \dfrac{3.4864}{2.689} = 1.2965$

and $\alpha^* = \bar{y} - \beta \bar{x} = (0.1161 - (1.2965) \cdot (1.5.37)$

$= 1.161 + 1.9631)$

$= 3.0442$

$\left[ \text{Here } \bar{Y} = \dfrac{\sum\limits_{i=1}^{n} Y_i}{n} = \dfrac{2.7773}{1} = 0.116, \text{ and } \bar{X} = \dfrac{\sum\limits_{i} X_i}{n} = \dfrac{4.0739}{2} = 1.5.37 \right.$

(ii) $R^2 = \dfrac{\beta \sum\limits_{i=1}^{n} X_i - \bar{X})^2}{\sum\limits_{i=1}^{n} Y_i - \bar{Y})^2} = \dfrac{\beta \sum\limits_{i=1}^{n} x_i^2}{\sum\limits_{i=1}^{n} (Y_i - \bar{Y})^2} = \dfrac{\beta^2 \sum\limits_{i=1}^{n} x_i^2}{\sum\limits_{i=1}^{n} y_i^2}$

$= \dfrac{(1.2965)^2 \times 2.6891}{4.8166} = 0.70$   $R^2 = 0.70$

(iii) $\sigma_u^2 = \dfrac{\sum e_i^2}{n-2} = \dfrac{\sum\limits_{i=1}^{n} y_i^2 - \beta^2 \sum x_i^2}{n-2}$

$= \dfrac{4.8566 - (1.2965)^2 \cdot 2.6891}{11-2} = \dfrac{4.8566 - 4.5201}{9} = \dfrac{3365}{9} = \ldots$

(iv) $Var(\beta) = \dfrac{\sigma_u^2}{\sum x_i^2} = \dfrac{0.0373}{2.6891} = 0.01387$

$SE(\beta) = \sqrt{0.01387} = 0.1177$

The regression results can now be written as follows

$$I^* = \alpha^* + \beta r^*$$

or, $\log I = 3.0812 - 1.12045 \log r \qquad \qquad R^2 = 0.76$
$\qquad \qquad (0.1177)$

or _____ where _____ antilog of 3.0812

The results show that the constant of the demand function ... ... means that the demand for investment is interest elastic. This result ... is consistent with economic theory.

## 2.19 Problem of Prediction / Forecasting Relating to a Two-variable Linear Regression Model

Today we do not differentiate between prediction and forecasting. We use the two of the other interchangeably. But these two terms are not identical. Prediction signifies an estimation of any event happening (in the past, present or future). On the other hand, forecasting is always associated with a time dimension, is the future estimation for some specific future duration or over a period of time. All forecasts are predictions, but not all predictions are forecasts, as when we use regression to explain the relationship between two variables. Forecast implies time series and where prediction does not. When we are interested to predict or forecast about future, we call the regression analysis as **historical regression**. With the help of regression analysis we can forecast about the future value on the basis of past and present information of the said variables ($X$ and $Y$). In the context of forecasting we may also distinguish between ex-ante forecast and ex-post forecast. Ex-ante forecast is a forecast that uses information available at the time of forecast, whereas ex-post forecast is a forecast that uses information beyond the time at which the forecast is made.

Let us define a classical linear regression model given by $Y = \alpha + \beta X + u$ for $i = 2, \ldots, n$ with the help of the pairs of observations $(X_1, Y_1), (X_2, Y_2), \ldots, (X_n, Y_n)$. We estimate the relationship by the method of least squares. The estimated relationship is $Y = \alpha + \beta X$. In case of time series data we write the regression equation as $Y_t = \alpha + \beta X_t + u_t$

Now for some value of $X$ (the independent variable), which is not in the sample, we may like to estimate the value of $Y$ (the dependent variable).

The process of finding the value of the dependent variable from the estimated relationship for the known value of the independent variable not in the sample is called "Prediction".

Let us suppose that $X_0$ is the value of the independent variable not in the sample and we have to predict the value of $Y$ when $X = X_0$. There are two types of prediction (i) Point prediction. (ii) Interval prediction

### 2.19.1 Point Prediction

When prediction is done in terms of a single value of the dependent variable, then it is called point prediction. We simply put $X = X_0$ in the estimated relationship and we get, $\hat{Y} = \hat{\alpha} + \hat{\beta} X_0 = Y_0$.

Now, the true value of the dependent variable $Y$ when prediction is made and it is given by

$$Y_0 = \alpha + \beta X_0 + u_0, \text{ where } u_0 \text{ is the corresponding value of the disturbance term}$$

So $\hat{Y}_0 = \hat{\alpha} + \hat{\beta} X_0 = E(\hat{Y}_0)$

$$\qquad = \hat{\alpha} + \hat{\beta} X_0 + E(u_0) = 0$$

Let us define the prediction error by $e_0 = Y_0 - \hat{Y}_0$ when we want to predict $Y_0$

and $e_0 = Y_0 - \hat{Y}_0$ is the prediction error when we want to predict $E_0$.

**Mean of the Predictor**

When we want to predict $Y_0$ then $e_0 = Y_0 - \hat{Y}_0$

Now $e_0 = Y_0 - \hat{Y}_0 = \alpha + \beta X_0 + u_0 - \hat{\alpha} - \hat{\beta} X_0$

or $e_0 = u_0 + (\alpha - \hat{\alpha}) + (\beta - \hat{\beta}) X_0$

or $E(e_0) = E(u_0) + (\alpha - E(\hat{\alpha})) + E(\beta - \hat{\beta}) X_0 = 0$

where $E(u_0) = 0$, $E(\hat{\alpha}) = \alpha$, $E(\hat{\beta}) = \beta$

Hence $E(Y_0 - \hat{Y}_0) = E(e_0) = 0$

or, $Y_0 - E(\hat{Y}_0) = 0$ or, $E(\hat{Y}_0) = Y_0$

So $\hat{Y}_0$ the OLS point predictor of $Y_0$ is unbiased

**Variance of the predictor**  Variance of the predictor is given by

$$\mathrm{Var}(\hat{Y}_0) = E\left[\hat{Y}_0 - E(\hat{Y}_0)\right]^2 = E\left[\hat{Y}_0 - Y_0\right]^2 \quad \left[\because E(\hat{Y}_0) = Y_0\right]$$

$$= E[e_0]^2 = E(e_0)^2 \quad [\because e_0 = Y_0 - \hat{Y}_0]$$

Now $E(e_0)^2 = E\left[u_0 + (\alpha - \hat{\alpha}) + (\beta - \hat{\beta}) X_0\right]^2$

$$= E[u_0^2 + (\alpha - \hat{\alpha})^2 + (\beta - \hat{\beta})^2 X_0^2 + 2(\alpha - \hat{\alpha}) u_0$$

$$+ 2(\beta - \hat{\beta}) X_0 u_0 + 2(\alpha - \hat{\alpha})(\beta - \hat{\beta}) X_0$$

$$= E u_0^2 + E(\alpha - \hat{\alpha})^2 + E(\beta - \hat{\beta})^2 X_0^2 + 2E(\alpha - \hat{\alpha}) u_0$$

$$+ 2X_0 E(\beta - \hat{\beta}) u_0 + 2X_0 E(\alpha - \hat{\alpha})(\beta - \hat{\beta})$$

$$= \sigma_u^2 + \mathrm{var}(\hat{\alpha}) + X_0^2 \mathrm{var}(\hat{\beta}) + 2\mathrm{cov}(\hat{\alpha}, u_0)$$

$$+ 2X_0 \mathrm{cov}(\hat{\beta}, u_0) + 2X_0 \mathrm{cov}(\hat{\alpha}, \hat{\beta})$$

We know that, $\mathrm{var}(\hat{\alpha}) = \sigma_u^2 \left[ \dfrac{1}{n} + \dfrac{\bar{X}^2}{\displaystyle\sum_{i=1}^{n} x_i^2} \right]$ and $\mathrm{cov}(\hat{\alpha}, \hat{\beta}) = \mathrm{cov}(\hat{\alpha}, \hat{\beta})$ as the

estimators of the parameters are independent of the disturbance term.

Now $\operatorname{cov}(\hat{\alpha}, \hat{\beta}) = E\left[(\hat{\alpha} - \alpha)(\hat{\beta} - \beta)\right]$

$$\operatorname{var}(\hat{\alpha}) = E(\hat{\alpha})^2 = \sigma_u^2 \left[\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^{n} x_i^2}\right]$$

$$\operatorname{var}(\hat{Y}_0) = E(e_0^2) = \sigma_u^2 \left(1 + \frac{1}{n}\right) + (X_0 - \bar{X})^2 \operatorname{var}(\hat{\beta})$$

Except $\operatorname{var}(\hat{\beta})$ all the terms are constant and positive

So, $\operatorname{var}(\hat{Y}_0)$ is minimum when $\operatorname{var}(\hat{\beta})$ is minimum

We know that, $\operatorname{var}(\hat{\beta})$ is minimum when $\hat{\beta}$ is the OLS estimator of $\beta$. Hence $\operatorname{var}(\hat{Y}_0)$ is minimum when $\hat{Y}_0$ is the OLS point-predictor of $Y_0$. This is the BLUE property of OLS point predictor. It means that OLS point predictor of $Y_0$ i.e. $\hat{Y}_0$ is the best linear unbiased predictor of $Y_0$.

We now consider another case where we want to make a point prediction of $E(Y_0)$.

Here prediction error is defined as $E(Y_0) - \hat{Y}_0 = e_0$. We know that $Y_0 = \alpha + \beta X_0 + u_0$

$E(Y_0) = \alpha + \beta X_0$ because $E(u_0) = 0$

$e_0 = \alpha + \beta X_0 - \hat{\alpha} - \hat{\beta} X_0$

$= (\alpha - \hat{\alpha}) + (\beta - \hat{\beta}) X_0 = -(\hat{\alpha} - \alpha) - (\hat{\beta} - \beta) X_0$

$e_0^2 = (\hat{\alpha} - \alpha)^2 + (\hat{\beta} - \beta)^2 X_0^2 + 2 X_0 (\hat{\alpha} - \alpha)(\hat{\beta} - \beta)$

or, $E(e_0^2) = E(\hat{\alpha} - \alpha)^2 + X_0^2 E(\hat{\beta} - \beta)^2 + 2 X_0 E(\hat{\alpha} - \alpha)(\hat{\beta} - \beta)$

$E(e_0^2) = E\left[\hat{Y}_0 - E(Y_0)\right]^2 = \operatorname{var}(\hat{\alpha}) + X_0^2 \operatorname{var}(\hat{\beta}) + 2 X_0 \operatorname{cov}(\hat{\alpha}, \hat{\beta})$

Since var $\ldots$ $= \ldots$ $\ldots$ and

$$\sum \qquad \sum$$

and cov $(\ldots, \beta) = \ldots$ var $\beta$.

$$\hat{Y} = \beta_0 + \beta_1 x_i, \quad \beta_1 = \ldots, \quad \text{or} \quad \bar{Y} = \hat{\beta}_0 \text{ and } x = \bar{Y} - x$$

$$\hat{Y} = \bar{x} \text{ is } \beta_0 \beta_1 (x_i) \quad \text{i.e. } \beta_0(x_i) \qquad + \text{var}\beta$$

Now putting the values of var $(\ldots)$ and cov $(\alpha, \beta)$ in the above expression we get

$$E\hat{Y}_0 \quad E\hat{Y}_{0i} \qquad \text{var}\hat{Y}_{0i} = \text{var} \, e_{0i} \qquad f(\hat{Y}, )$$

Now var $e_0 = E\,\hat{Y}_{0i} + \hat{Y}_0$

$$\frac{\sigma_u^2}{n} + \hat{x}_0 \, \text{var}\beta \qquad \hat{x}_0 \, \text{var}\beta) \qquad \bar{x}_0 \hat{x} \, \text{var}\beta \qquad \frac{\sigma_u}{n} + \hat{x}_0 \quad \hat{x}_0 \quad \hat{x}_0 \, \text{var}\beta$$

var $(e_0)$ is maximum when var$\beta$ is minimum. Now var$\beta$ is minimum where

the OLS estimator of $\beta$. So var$(e_0)$ is minimum when $\hat{Y}_0$ is the OLS point prediction

$E(Y_0)$. This is the BLUE property of the OLS point prediction $\hat{Y}_0$.

## 2.19.2 Test of Significance of Predictor and Interval Prediction

**Case 1** We want to test the null hypothesis $H_0: Y_0 = $ A some specified value

against the alternative hypothesis $H: Y_0 \neq A$ or $H: Y_0 < A$ or $H: Y_0 > A$. We use

$\hat{Y}_0$ as the appropriate statistic of $Y_0$ because $\hat{Y}_0$ is the BLUE predictor of $Y_0$. Now

$\hat{Y}_0 = \alpha + \beta x_0$. Since $\hat{Y}_0$ is a linear function of $\alpha$ and $\beta$, $\hat{\alpha}$ and $\hat{\beta}$ are normally

distributed. So $\hat{Y}_0$ is also normally distributed.

Since $E(e_0) = E\left[\hat{Y}_0 - Y_0\right] = 0, \quad E(\hat{Y}_0) = Y_0$

and var$(\hat{Y}_0) = E\left[\hat{Y}_0 - Y_0\right]^2 = F(e_0)$.

$$= \sigma_u^2 - \text{var}(\alpha) + \text{var}\beta)\, Y_0^2 - 2\hat{x}_0 \, \text{cov}(\alpha, \beta)$$

$$= \sigma_u^2 + \sigma_u^2 \left[\frac{1}{n} + \frac{\hat{x}}{\sum_{i=1}^{n} x_i}\right] \quad \hat{x}_0 \frac{\sigma_u}{\sum x_i^2} \quad 2\hat{x}_0 \bar{Y} \frac{\sigma_u}{\sum x_i^2} \quad \text{where cov}(\alpha, \beta) \quad \bar{x}\,\text{var}$$

and var$(\beta) = \dfrac{\sigma_u}{\sum\limits_{i=1}^{n} x_i^2}$

$$\operatorname{var}(\hat{Y}_0) = \sigma_\epsilon^2 \left[ 1 + \frac{1}{n} + \frac{1}{\sum_{i} c_i^2} (\bar{X} - X_0)^2 \right]$$

This means that $\hat{Y}_0$ is normally distributed with mean $Y_0$ and variance

$$\sigma_\epsilon^2 \left[ \ldots \right] \frac{1}{\sum_{i} c_i^2}$$

So if $\sigma_\epsilon$ is unknown it is to be replaced by its unbiased

estimator $\sum_{i} e_i^2 / (n-2)$. The appropriate test statistic will be given by

$$\frac{\hat{Y}_0 - Y_0}{\sqrt{\frac{\sum_{i} e_i^2}{n-2} \left[ \ldots + \frac{(\bar{X} - Y_0)^2}{\sum_{i=1}^{n} e_i^2} \right]}} = t_{n-2}$$

it follows a '$t$' distribution with $(n-2)$ degrees of freedom.

**Nature of the test** If the alternative hypothesis is $H : Y_0 \ne A$ then the null hypothesis will the accepted at 5% level of significance if $t_{0.025, n-2} \le t \le t_{0.025, n-2}$ and will be rejected otherwise.

If the alternative hypothesis is $H_1 : Y_0 > A$ then $H_0 : Y_0 = A$ will be accepted at 5% level of significance if $t$ (observed) $\le t_{0.05, n-2}$ (table) and will be rejected otherwise.

If the alternative hypothesis is $H_1 : Y_0 < A$, then $H_0 : Y_0 = A$ will be accepted at 5% level of significance if

$t$ (observed) $\ge -t_{0.05, n-2}$ (table), and will be rejected otherwise

The rejection of the null hypothesis on the basis of the sample data implies the significance of $Y_0$.

It should be noted that the same procedure can be used for 1% level of significance

It can be seen that $100(1-\alpha)\%$ confidence interval of $Y_0$ would be

$$Y_0 + Y_0 \pm \sqrt{\frac{\sum e_i^2}{n-2} \cdot \left[\cdots\right]} \cdot t_{\alpha/2}$$    where $\alpha$ usually takes the value as $.05$ or

$.01$.

**Case 2.** We want to test the null hypothesis $H_0: E(Y_0) = A$ against the alternative
$H: E(Y_0) \neq A$ or $H: E(Y_0) > A$ or $H: E(Y_0) < A$

Here we take $\hat{Y}_0$ as the estimate of $E(Y_0)$ because $\hat{Y}_0$ is the BLUE predictor of $E(Y_0)$

Here also $\hat{Y}_0$ is normally distributed with mean $E(Y_0)$ and variance

$$\text{var } \hat{Y}_0 = \frac{\sigma_u^2}{n} + (X_0 - \bar{X})^2 \text{var}(\hat{\beta})$$

$$= \frac{\sigma_u^2}{n} + X_0^2 \cdot \sigma_u^2 \frac{1}{\sum_{i=1}^{n} x_i^2} = \sigma_u^2 \left[ \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum_{i=1}^{n} x_i^2} \right]$$

So, $\hat{Y}_0$ is normally distributed with mean $E(Y_0)$ and variance $\sigma_u^2 \left[ \dfrac{1}{n} + \dfrac{(X_0 - \bar{X})^2}{\sum_{i=1}^{n} x_i^2} \right]$

If $\sigma_u^2$ is unknown, it is to be replaced by its unbiased estimator $\dfrac{\sum_{i=1}^{n} e_i^2}{n-2}$. The

appropriate test statistic under $H_0: E(Y_0) = A$ would be

$$\frac{\hat{Y}_0 - E(Y_0)}{\sqrt{\dfrac{\sum_{i=1}^{n} e_i^2}{n-2} \left[ \dfrac{1}{n} + \dfrac{(X_0 - \bar{X})^2}{\sum_{i=1}^{n} x_i^2} \right]}} = t_{n-2}$$

which follows a $t$-distribution with $(n-2)$ degrees of freedom.

**Nature of the test**

If $H_0: E(Y_0) = A$ is tested against the alternative $H: E(Y_0) \neq A$, then $E(Y_0) = A$ will be accepted at $5\%$ level of significance if $t_{0.025, n-2} \leq$ observed $< t_{0.025, n-2}$ and will be rejected otherwise.

I the alternative hypothesis is $H$ ... then ... is ... we accept at the level of significance if observed ... value ... and ... otherwise.

the alternative hypothesis is $H$ ... then ... is ... we accept at the level of significance if observed ... value ... and ... otherwise. The same test procedure is applicable at $H$ level of significance.

It can be seen that $90(1 - \alpha)\%$ confidence interval of ... will be

$$t_0 + \alpha = \sqrt{\frac{\sum e^2}{n - 2} \left( 1 + \frac{1}{n} \right) \cdot \frac{1}{\sum}} \qquad \text{where } n \text{ usually takes the value } 0.05 \text{ or}$$

**Example 3.11.** Following example 2.3.

(i) find out the point predictor of $Y$ when $X = 10$

(ii) It is assumed that when $X = 10, Y = 65$. Do you think that it is justified?

(iii) It is claimed that when $X = 10, E(Y) = 55$. Do you think that the claim is justified?

**Solution**

(i) We have to find out the point predictor of $Y$, when $X = 10$. We know that the point predictor of $Y$, is.

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i \qquad \qquad \text{(where } \hat{\alpha} = 3.505, \ \hat{\beta} = 0.494 \text{ see Ex 2.3)}$$

$$= -3.505 + 0.494 \times 10$$

$$= 3.505 + 4.94 = .435$$

$$\hat{Y} = .435$$

So, point predictor of $Y$, is $\hat{Y} = 1.435$ when $X = 10$

(ii) We have to examine whether $Y = 165$ when $X = 10$ is justified or not. We have to test the null hypothesis $H_0 : Y_0 = 165$, against the alternative $H : Y_0 \neq .65$

The appropriate test statistic is given by, 
$$\frac{Y_0 - \hat{Y}_0}{\sqrt{\frac{\sum_{i=1}^{n} e_i^2}{n-2} \left[ 1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{\sum_{i=1}^{n} x_i^2} \right]}} \qquad t_{n-2}$$

Since $\hat{\alpha} = 3.505$, $\hat{\beta} = 0.494$ and $X_0 = 10$,

then $Y_0 = \hat{\alpha} + \hat{\beta} X_0 = 3.505 + 0.494 \times 10 = 1.435$

$\hat{Y}_0 - Y_0 = .435 - 165.000 = -163.565$

$$\hat{\beta}_1 = \frac{\sum c_i}{n}$$

$$\text{Since } \bar{c} = \frac{\sum c_i}{n} = \ldots$$

$$\frac{\ldots}{k} = \ldots$$

$$= \ldots$$

$$= \sqrt{90(\ldots)} \ldots \ldots$$

$$= \sqrt{900 \ldots (0.5 + 0.02 \ldots)}$$

$$= \sqrt{900 \ldots 1577} = \sqrt{900 \, 76.67}$$

$$= 4.69$$

$$\text{and } \sum_{i=1}^{n} c_i = 2.54$$

Now $t$ (observed) $=$

$$\frac{Y_0 - Y_0}{\sqrt{\dfrac{\sum \hat{e}_i^2}{n-2}\left[\dfrac{1}{n} + \dfrac{(\bar{X} - X)^2}{\sum_{i=1}^{n} x_i^2}\right]}}$$

$$= \frac{-61.565}{4.69} = -115.436$$

$$t = -5.438$$

Here we see that $t$ (observed) $= -115.436$ which does not lie in the interval $-t_{0.025}$, 18 and $t_{0.025}$, 18, i.e. in the interval $-2.101$ and $2.101$, and hence the null hypothesis $H_0 : Y_0 = 165$ is rejected for the given sample at 5% level of significance. So, $Y_0 = 165$ is not justified when $X_0 = 10$.

(b) We have to examine whether $E(Y_0) = 155$ when $X_0 = 10$.

We have to test the null hypothesis $H_0 : E(Y_0) = 155$, against the alternative $H_1 : E(Y_0) \neq 155$.

The statement $E(Y_0) \neq 155$, when $X_0 = 10$ will be justified if the null hypothesis $H_0 : E(Y_0) = 155$ is rejected.

The appropriate test statistic would be

$$t = \frac{Y_0 - E(Y_0)}{\sqrt{\dfrac{\sum \hat{e}_i^2}{n-2}\left[\dfrac{1}{n} + \dfrac{(X_0 - \bar{X})^2}{\sum_{i=1}^{n} x_i^2}\right]}} \sim t_{n-2}$$

Here , $n\bar{x}$ = 419, $\sum Y_0$ = 155

$\sum$ ...

, $x_0$ = 10, $\bar{x}$ = 0.31, $\sum x_i^2$ = 215.4

$$\sqrt{\frac{\sum Y_i}{n-2}\left[\frac{1}{n}+\frac{x_0^2}{\sum x_i^2}\right]}$$

$$\sqrt{\frac{408}{20}+\frac{(0.91)}{215.4}}$$

$= \sqrt{1.908[0.05 + 0.06221]}$, $= \sqrt{1.908 \times 0.15221}$ ... 0.1156

Now $t$ (observed) $= \dfrac{\bar{Y}_0 - E(Y_0)}{\sqrt{\dfrac{\sum x_i^2}{n-2}\left[\dfrac{1}{n}+\dfrac{(x_0-\bar{x})^2}{\sum\limits_{i=1}^{n} x_i^2}\right]}}$    $\dfrac{-53.465}{0.156}$ , $-446.581$

$t$ (observed) $= -486.581$

Now on the basis of the given sample the null hypothesis $H_0$ : $E(Y_0) = 155$ will be accepted at 5% level of significance if $-t_{0.025,n-2} < t < t_{0.025,n-2}$ and will be rejected otherwise.

Here $t_{0.025,n-2} = t_{0.025,18} = 2.101$

So, the observed $t = -486.581$ does not lie in the interval $-2.101$ and $2.102$ and hence the null hypothesis will be rejected. So, $E(Y_0) = 155$ when $x_0 = 10$ is not justified.

**Example 2.12.** Consider the following regression model $Y = \alpha + \beta x + u$ where $u$ is normally distributed with mean zero and variance $\sigma_u^2$ (unknown). We have the following data

| X | 2 | 3 | 1 | 5 | 9 |
|---|---|---|---|---|---|
| Y | 4 | 7 | 3 | 9 | 7 |

(i) Estimate $\alpha$ and $\beta$.

(ii) Test whether $\alpha$ and $\beta$ are significant or not at 5% level of significance.

(iii) Calculate $R^2$

... ...

**Solution:** ... ... estimate the regression parameters $\alpha$ and $\beta$ ... ... and ... ... ... ... ... ... ... ... ...

$$\sum \qquad \sum \qquad \text{...} \qquad \beta\bar{X}$$

### Calculations for $\alpha$, $\beta$ and $R^2$

| | | | | | $-X$   $\beta$ | $Y$   $Y$ | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 6 | 8 | 4 | 2 | ? | 4 | 4 | 4 |
| | 5 | 7 | 9 | | | 1 | 1 | 2 |
| 5 | 7 | 7 | 1 | 3 | | 1 | 0 | |
| | 6 | 40 | 24 | 1 | 1 | 1 | 14 | 1 |
| | 7 | 63 | 30 | 5 | 1 | 25 | 1 | 5 |
| $\sum X$ | $\sum Y$ | $\sum XY$ | $\sum X^2$ | $\sum u$ | $\sum Y$ | $\sum \hat{Y}$ | $\sum e$ | $\sum$ |
| $=20$ | $=30$ | $=40$ | $=20$ | $=0$ | $=0$ | $=40$ | $=24$ | $=24$ |

Here $\bar{X} = \sum \frac{X}{n} = \frac{20}{5} = 4$ and $\bar{Y} = \sum \frac{Y}{n} = \frac{30}{5} = 6$

Then $\beta = \sum XY / \sum X^2 = \frac{20}{40} = 0.5$ and $\alpha = \bar{Y} - \beta\bar{X} = 6 - 0.5 \times 4 = 6 - 2 = 4$

$\alpha = 4 \quad \beta = 0.5$

Thus the OLS estimators of $\alpha$ and $\beta$ are $\hat{\alpha} = 4$ and $\hat{\beta} = 0.5$

The estimated regression line is $\hat{Y} = \alpha + \beta X$ or $\hat{Y} = 4 + 0.5X$

Now we have to find out $Var(\alpha)$ and $Var(\beta)$. We know that

$$Var(\alpha) = \sigma_u^2 \frac{\sum\limits_{i=1}^{n} X^2}{n\sum\limits_{i=1}^{n} x_i^2} \quad \text{and} \quad Var(\beta) = \sigma_u^2 \Big/ \sum\limits_{i=1}^{n} x_i^2$$

But here $\sigma_u^2$ is not known and hence it is to be replaced by its unbiased estimator

$$\hat{\sigma}_u^2 = \sum\limits_{i=1}^{n} e_i^2 \Big/ (n-2)$$

$$\hat\alpha = \bar{y} - \hat\beta \bar{x}$$

For $\hat\alpha = 2.802$ and $SE(\hat\alpha) = \sqrt{var(\hat\alpha)} = \sqrt{2.802} = 1.674$

Similarly $var(\hat\beta) = \sigma_u^2 \sum_{i=1}^{n} x_i^2 = \dfrac{4.67}{40} = 0.11675$

For $var(\hat\beta) = 0.11675$ and $SE(\hat\beta) = \sqrt{var(\hat\beta)} = \sqrt{0.11675} = 0.3416$

(ii) **Test for $\hat\alpha$ and $\hat\beta$ :**

a **Test for $\beta$** : We have to test the null hypothesis $H_0 : \beta = 0$ against the alternative $H : \beta \neq 0$. The appropriate test statistic would be

$$t = \frac{\hat\beta}{SE(\hat\beta)} \sim t_{n-2}$$

The null hypothesis will be accepted at 5% level of significance if $t_{0.025,n-2} > t > t_{0.025,n-2}$ and will be rejected otherwise.
Here we see that,

$$t = \frac{\hat\beta}{SE(\hat\beta)} = \frac{\hat\beta}{\sqrt{\left[ \sum_{i=1}^{n} e_i^2 / (n-2) \right] / \sum_{i=1}^{n} x_i^2}} = \frac{0.5}{0.3416} = 1.4637$$

$t$ (observed) $= 1.4637$

But from table value $t_{0.025,n-2} = t_{0.025,5-2} = t_{0.025,3} = 3.182$

Here we see that the $t$ (observed) $= 1.4637$ lies in the interval $-3.182$ and $3.182$ and hence the null hypothesis is accepted at 5% level of significance.

So, $\beta$ is insignificant at 5% level — significant only when the null hypothesis is rejected.

**b) Test for $\bar{\alpha}$.** If we test the null hypothesis $H_0 : \alpha = 0$ against the alternative ... then the null hypothesis will be accepted at 5% level of significance ... and will be rejected otherwise ...

From above values ... $t = \frac{4}{1.674} = 2.3895$

and consequently $\frac{t}{\sqrt{...}}$ ...

$$\sum \hat{\varepsilon} \quad \sum 1$$

Thus we see that $t$ observed $= 2.3895$ lies in the interval $-t_{0.5, 11}$ and $t_{0.5, 11} = -2.82$ and $2.82$, and hence the null hypothesis is accepted at 5% level of significance. This means that it is also insignificant at 5% level of significance.

Now we have to calculate the value of $R^2$.

Since we know that $R^2 = \frac{ESS}{TSS} = \frac{\hat{\beta} \sum\limits^{n} x_i^2}{\sum\limits^{n} y_i^2} = \frac{(0.5)^2 \cdot 40}{24} = \frac{10}{24} = 0.4167 \approx 0.42$

$R = 0.42 = \frac{42}{100} = \frac{\text{Explained variation}}{\text{Total variation}}$

This suggests that 42 percent of the variations in the sample observations of $Y$ can be attributed to the variations of the fitted value of $Y$ ($\hat{Y}$) or we can say that our regression line fits the given data not very well.

From the above results we can write our regression results as follows

$$Y = 4 - 0.5X \quad R^2 = 0.42$$
$$SE \quad (1.674) \ (0.3416)$$

Alternatively, $\quad Y = 4 + 0.5X \quad R^2 = 0.42$

$t \text{ ratios} \quad (2.3895) \ (1.4637)$

(v) We have to find out the point predictor of $Y$ when $X = 10$

The point predictor of $Y$ is given by $\hat{Y} = \alpha + \beta X$

$= 4 + 0.5 \times 10 = 4 + 5 = 9$

Point predictor $Y = 9$ when $X = 10$

**Example 2.13.** Following the data given in Example 2 .

(i) estimate the regression parameters (assuming a linear regression equation of the

form $Y = \alpha + \beta X_i + u_i$ where $u_i \sim N(0, \sigma_u^2)$ (unknown)

(b) calculate $R$

(ii) ... obtain the ... ... expenditures are ... ... in $Y$ ...

(c) find out RSS ... ... ... when $Y = 4.44$

... and $0.74$ ... find a value of the estimated predicted value of $Y$ when $Y = 4.44$

**Solution**

Table for calculation

| Month | 1 | 2 | x | y | | | $xy$ | $\hat{Y}$ | $\hat{e}$ |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 2 | 6 | | 4 | | | 0.8 | 0.44 |
| 2 | 3 | 4 | 1.6 | 1.6 | | | | 0.6 | 0.36 |
| 3 | 1 | 2 | 0 | 2.6 | 5 | 0 | 6.6 | 6 | 6.76 |
| 4 | 4 | 6 | 1 | 4 | 6 | | | | 6.64 |
| 5 | 4 | 8 | 2 | 1.6 | 6.1 | 4 | 2.4 | | 0 |
| Total | 2.3 | 2.7 | 2.1 | 2 | | 1 | 2 | | |
| | 15 | 23 | 1 | 0 | | | 2.0 | | 4.4 |

Here $n = 5$ as we have data for five months.

Now $\bar{X} = \dfrac{\sum x}{n} = \dfrac{15}{5} = 3$, $\bar{Y} = \dfrac{\sum Y}{n} = \dfrac{23}{5} = 4.6$

(i) We have to estimate the regression parameters $\alpha$ and $\beta$. Let $\hat{\alpha}$ and $\hat{\beta}$ be the OLS estimators (predictors here) of $\alpha$ and $\beta$.

We know that $\hat{\beta} = \dfrac{\sum x_i y_i}{\sum x_i^2}$ and $\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}$

$\hat{\beta} = \dfrac{12}{10} = 1.2$ and $\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X} = 4.6 - 1.2 \cdot 3 = 4.6 - 3.6 = 1$

Therefore the estimated (predicted) sample regression equation is $\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i = 1.0 + 1.2X_i$

In the table $\hat{e} = Y_i - \hat{Y}_i$ can be obtained for different values of $Y$. Since $\hat{e} = Y - 1.0 - 1.2X_i$.

When $X_i = 1$ and $Y_i = 3$, $\hat{e}_i = 3 - 1.0 - 1.2 \times 1 = 3.0 - 2.2 = 0.8$

$X_i = 2$ and $Y_i = 4$, $\hat{e}_i = 4 - 1.0 - 1.2 \times 2 = 4 - 1 - 2.4 = 0.6$

and in this way other $\hat{e}_i$ values are calculated.

(ii) We know that $R^2 = \dfrac{ESS}{TSS} = \dfrac{\hat{\beta}^2 \sum\limits_{i=1}^{n} x_i^2}{\sum\limits_{i=1}^{n} x_i^2} = \dfrac{(1.2)^2 \times 10}{23.20} = \dfrac{14.4}{23.20} = 0.620$

Since $\sum\limits_{i=1}^{n} y_i^2 = \hat{\beta}^2 \sum\limits_{i=1}^{n} x_i^2 + \sum\limits_{i=1}^{n} \hat{e}_i^2$ i.e. TSS = ESS + RSS $\qquad R^2 = 0.620$

...we predicted regression ...... in ...... by $Y$ .... $\hat{Y}$ ....

.... where ...... $\hat{Y}_0$ ...... then ......

...... ......

This ..... by .... of ...... expenditures are in respect to ₹ 600 .... sales $=$ ...... between ₹ .....

.... We have to find out 95% confidence interval of the predicted value of $\hat{Y}$ when $X = 6$ now

We ...... that at 95% ...... $(\alpha = 0.05)$ here confidence interval of the ...... (predicted) value of the ...... product is $M$

$$\hat{Y}_0 \pm t_{\alpha,...} \sqrt{..... \left[1 + \frac{1}{n} + \frac{(.....)^2}{\sum x_i^2}\right]}$$

Here we have $\hat{Y}_0 = 4.5$ when $\bar{X} = X_0 = 6$

and $\sum_{i=1}^{n} e^2 = 8.8$, $n = 5$, $\sum_{i=1}^{n} x_i^2 = 10$, $\bar{X} = 3$, $\dfrac{\sum_{i=1}^{n} e_i^2}{(n-2)} = \dfrac{8.8}{3} = 2.93$

and $t_{n-2,\alpha} = t_{(3,0.05)} = 3.182$ [From table value]

So, 95% confidence interval of the predicted value of sales revenue corresponding to advertising expenditure of ₹ 600 would be

$$= 20 \pm 3.82 \sqrt{2.93\left[1 + \frac{1}{5} + \frac{(6-3)^2}{10}\right]}$$

or, $\, 20 \pm 7.891$ or $0.309$ and $16.09$

$\Rightarrow$ ₹ 309 and ₹ 16091

Thus 95% confidence interval of predicted sales revenue $\hat{Y}_0$ corresponding to advertising expenditure of ₹ 600 would be ₹ 309 and ₹ 16091

(vi) 95% $= (1 - \alpha)\%$ (a% when $\alpha = 0.05$) confidence interval of expected sales revenue $E(\hat{Y}_0)$, when advertising expenditures are ₹ 600 would be

$$\hat{Y}_0 \pm t_{\alpha/2,n} \sqrt{\frac{\sum_{i=1}^{n} e_i^2}{n-2}\left[\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum x_i^2}\right]}$$

[Here $\alpha = 0.05$, $t_{\alpha/2,n} = t_{(3,0.05)} = 3.182$ [from table value]

$n = 5$, $X_0 = 6$, $\sum x_i^2 = 10$, $\bar{X} = 3$, $\dfrac{\sum_{i=1}^{n} e_i^2}{(n-2)} = 2.93$

The prediction is still given by

$\hat{Y}_0 = 10 + 1.2X$, so that when $X = 1 = 6$, $\hat{Y}_0 = 10 + 1.2 \times 6 = 8.20$)

or, $8.20 \pm 3.182 \times 1.795$
or, $8.20 \pm 5.72$ or, 2.48 and 13.92

95% confidence interval for the average sales $E(Y_0)$ corresponding to advertising expenditures ₹ 600 would be ₹ 2480 and ₹ 13920.

It should be noted that this confidence interval is narrower than the one we obtained for $\hat{Y}_0$.

**Example 2.14** The following table (Table 2.6) gives data on the level of education (number of years of schooling), the mean hourly wages earned by the people at each level of education and the number of people at the stated level of education.

**Table 2.6. Mean Hourly wage by Education**

| Years of schooling (X) | Mean hourly wage in $ (Y) | Number of people |
|---|---|---|
| 6 | 4.4567 | 3 |
| 7 | 4.7700 | 5 |
| 8 | 5.9763 | 15 |
| 9 | 7.3317 | 12 |
| 10 | 7.3182 | 17 |
| 11 | 6.5844 | 27 |
| 12 | 7.8182 | 218 |
| 13 | 7.8384 | 37 |
| 14 | 11.0223 | 56 |
| 15 | 10.6738 | 13 |
| 16 | 10.8361 | 70 |
| 17 | 13.6150 | 24 |
| 18 | 13.5310 | 31 |
| Total | | 528 |

(i) Assuming a linear regression line of the form $Y_i = \alpha + \beta X_i + u_i$ [$u_i \sim N(0, \sigma_u^2)$], find the OLS estimators of $\alpha$ and $\beta$.

(ii) Find var $(\hat{\alpha})$ and var $(\hat{\beta})$.

(iii) Find $SE(\hat{\alpha})$ and $SE(\hat{\beta})$.

(iv) Find $R^2$.

(v) Find $\sum_{i=1}^{n} e_i^2$.

(vi) Predict/Forecast about the mean hourly wage when the level of education (years of schooling) is 20.

### Calculations for the Regression

| | | | | | | | |
|---|---|---|---|---|---|---|---|

(table illegible)

Total $\sum Y$  $\sum x$  $\sum x_i$  $\sum$  $\sum$  $\sum$

**Note** $n = ?$  $\bar{Y} = \sum Y / n$  $\dfrac{2712}{13} = 8.6742$

and $\bar{X} = \sum X / n = \dfrac{56}{13} = 12$  $\sum x_i = (6)^2 + (7) + \cdots + (8) = 2.54$

$$\beta = \dfrac{\sum x_i y_i}{\sum x_i^2} = \dfrac{7856}{1820} = 0.7240967$$

and $\alpha = \bar{Y} - \beta\bar{X} = 8.6747 - 0.7240967 \times 12 = -0.01445$

The estimated regression equation is

$$\hat{Y} = \alpha + \beta X \quad \text{or} \quad Y = -0.01445 - 0.7240967 X_i$$

$$e_i = Y - \hat{Y}_i = Y - 0.01445 - 0.7240967 X$$

Now different values of $e_i$ can be obtained by taking different pairs of values and $Y_i$

When ... when values of ... are obtained

$$N = \sum y \ldots$$      ...

$$\text{and} \sum \ldots$$

(ii) OLS estimates of $\alpha$ and $\beta$ would be

$$\beta = \frac{\sum x_i y_i}{\sum x_i^2} = \frac{1 \cdot 7850}{82} = 0 \cdot 240967$$

and $\alpha = \bar{y} - \beta \bar{x} = 8.6742 - 0.240967 \ldots = 0.01445$

$$\alpha = 0.0144 \text{ and } \hat{\beta} = 0.7240$$

(iii) We have to calculate the values of $\text{var}(\alpha)$ and $\text{var}(\beta)$

$$\sigma_u^2 \sum X$$

We know that $\text{var}(\alpha) = \dfrac{\sigma_u^2 \sum\limits_{i=1}^{n} x^2}{n \sum\limits_{i=1}^{n} x^2}$. Here $\sigma_u^2$ is unknown and hence it is replaced

by its unbiased estimator $\sigma_u^2 = \sum\limits_{i=1}^{n} e_i^2 / (n-2) = 0.8936$

$$\text{var}(\alpha) = \frac{\hat{\sigma}_u^2 \sum\limits_{i=1}^{n} Y^2}{n \sum\limits_{i=1}^{n} e_i^2} = \frac{0.8936 \times 7054}{13 \times 182} = \frac{1835.4544}{2366} = 0.7757$$

$$\text{var}(\hat{\alpha}) = 0.7757$$

Again, $\text{var}(\hat{\beta}) = \dfrac{\sigma_u^2}{\sum\limits_{i=1}^{n} x_i^2}$. Here $\sigma_u^2$ is unknown and hence it is replaced by its unbiased

estimator $\sigma_u^2 = \sum\limits_{i=1}^{n} e^2 / (n-2) = 0.8936$

$$\text{var}(\hat{\beta}) = \frac{\hat{\sigma}_u^2}{\sum\limits_{i=1}^{n} x_i^2} = \frac{0.8936}{182} = 0.004910 \qquad \text{var}(\beta) = 0.004910$$

$$\sum_{i} \quad (0.7240 \ldots \quad \sum e^2 = 9.830.7$$

$$= 94.6000 + 9.830.7 \quad 05.230.7$$

Note $\sum_{i=1} Y_i - \sum_{i=1} Y_i \ldots = (0.79\,406)^2 - (0.853\,117)^2 - \ldots + 0.0469\,TR_T$

$$= 9.83017 \qquad \sum_{i=1}^{n} e_i^2 = 9.83017$$

(iv) We have to forecast/predict about the mean hourly wage rate when the level education (years of schooling) is 20.

Since the estimated regression equation is $\hat{Y}_i = 0.0144 + 0.7240\,X_i$,

The point predictor of $Y$ is given by $\hat{Y} = \hat{a} + \hat{\beta}X$

When $X = X_0,\ \hat{Y}_0 = \hat{a} + \hat{\beta}X_0$

So, when the level of education is $X_0 = 20$, the mean hourly wage rate would be

$\hat{Y}_0 = 0.0144 + 0.7240 \times 20 = 14.4656$

So, mean hourly wage rate would be \$14.4656 when years of schooling increase to 20.

(v) We have to construct 95% confidence interval for the point predictor of hourly wage rate $(Y_0)$ when the level of education becomes $X_0 = 20$. This confidence interval would be

$$\hat{Y}_0 \pm t_{0.025,\,n-2} \sqrt{\frac{\sum_{i=1}^{n} e_i^2}{n-2}\left[1 + \frac{1}{n} + \frac{(\bar{X} - X_0)^2}{\sum_{i=1}^{n} x_i^2}\right]}$$

When $X = X_0 = 20$, $Y_0 = \hat{Y}_0 =$ ... 0.0144 + 0.7240 × 20 = 14.4656

... (Table ... )

Again ... $\sum_{i=1}$ ... $\sum_{i=1}$ ...

... 95% confidence interval of ... 

$$Y_0 \pm t_{0.025,n-2} \cdot s \sqrt{ \left[ 1 + \frac{1}{n} + \frac{(\bar{X} - X_0)^2}{\sum x_i} \right] }$$

or $14.4656 \pm 2.201 \sqrt{0.8936 \left[ 1 + \frac{1}{13} + \frac{(12 - 20)}{182} \right]}$

or $4.4656 \pm 2.201 \sqrt{0.8936 \left[ 1 + \frac{1}{13} + \frac{64}{182} \right]}$

or $4.4656 \pm 2.20 \sqrt{0.8936[1 + 0.07692 + 0.35164]}$

or $4.4656 \pm 2.20 \sqrt{0.8936 \times 1.42856}$ or $14.4656 \pm 2.201 \sqrt{1.27656}$

or $4.4656 \pm 2.201 \times 1.12984$ or, $14.4656 \pm 2.4868$

i.e., 9788 and 16.9524

95% confidence interval of hourly wage rate would be $ 9788 and $ 6.9524 when the level of education (years of schooling) is 20.

(v.) We have to construct 95% confidence interval of expected mean hourly wage rate when the level of education is $X_0 = 20$ (years of schooling).

When $X = X_0 = 20$, $E(Y|X_0 = 20)$ can be obtained as $\hat{Y}_0 = \alpha + \beta X_0 = 0.44 + 0.7240 \times 20 = 14.4656$

Thus 95% confidence interval of $E(Y|X_0)$ when $X_0 = 20$ would be

$$Y_0 \pm t_{0.025,n-2} \sqrt{ \frac{\sum_{i=1}^{n} e_i^2}{n-2} \left[ \frac{1}{n} + \frac{(\bar{X} - X_0)^2}{\sum_{i=1}^{n} x_i^2} \right] }$$

Here $\hat{Y}_0 = 14.4656$, $n = 13$, $t_{0.025,n-2} = t_{0.025,11} = 2.201$ (Table value)

$\bar{X} = 12$, $X_0 = 20$ and $\sum_{i=1}^{n} x_i^2 = 182$, $\sum_{i=1}^{n} e_i^2 / (n-2) = 0.8936$

$$\text{Now, } \hat{Y}_0 = \bar{Y}_0 \pm t \sqrt{\frac{\sum \hat{e}_i^2}{n-2}\left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum x_i^2}\right]}$$

$$= 14.4656 \pm 2.20 \sqrt{0.8936}\left[\frac{1}{11} + \frac{(1)^2}{78}\right]$$

or $14.4656 \pm 2.20 \sqrt{0.8936} \times 0.1243$ or $14.4656 \pm 2.20 \sqrt{0} \times 0.2470$

or $14.4656 \pm 2.20 \times 0.9378$ or $14.4656 \pm 1.7020$

c., 13, 030 and 18276

Thus, 95% confidence interval of expected mean hourly wage rate corresponding to the level of education (year of schooling) $X_0 = 20$ would be $ 13.030 and $ 16.2..

**Example 2.49.** The following table (Table 2.7) shows consumption expenditure and income in billions of $ of a country over the period 2007 to 2018.

**Table 2.7 Consumption expenditure and income of a country (in billions of $)**

| Year | Consumption expenditure $C_t$ | Income $Y_t$ |
|---|---|---|
| 2007 | 282.3 | 371 |
| 2008 | 291.1 | 394.7 |
| 2009 | 323.5 | 4.18 |
| 2010 | 3.6.6 | 430.8 |
| 2011 | 316.6 | 456.7 |
| 2012 | 366 | 486.5 |
| 2013 | 66 | 4.55 |
| 2014 | 386.0 | 55.4 |
| 2015 | 544 | 6.1 |
| 2016 | 3003.5 | 648.4 |
| 2017 | 3003.5 | 2.1 |
| 2018 | 4439.0 | |

From the data given in the table we have the following results

$$\hat{Y}_t = 1.76 + 0.73 Y, \quad R^2 = 0.998, \sigma_e^2 = \sum e_t^2 / n - 2 = 284.6$$

$(5.79) \quad (0.0.2)$

$\bar{X} = 498, n = 2$

$$\sum_{t=1}^{n} (X_t - \bar{X})^2 = \sum_{t=1}^{n} x_t^2 = 51.482, \quad C = \bar{Y}_0 = 273.904$$

$t_0$ = Income in the year 2025 = 5850 billion

(i) Forecast about consumption expenditure of the country for the year 2025 income in that year increases to 5854 billion

(ii) Construct 95% confidence interval of the consumption expenditure Predicted forecasted for the year 2025

Solution

$\hat{Y} = 31.78 + 0.71 X_t$, or when $X_t = 810 + I_t$

$= 31.78 + 0.71 \times 810 = 631.33$ billion

...

$$635 + 6 \quad \sqrt{265.61\left[1 + \frac{1}{12} + \frac{123.904}{151,482}\right]}$$

or $\quad 635 + 2.228 \times \sqrt{285.6} \times 18.4$

i.e. $\quad 635 + 2.228 \times 514 \times 8885$

or $\quad 635 + 2.228 \times 2 \quad 097$

i.e. $\quad 635 + 5.4884$

or $\quad 583.5$ to 686.4884

Thus 95% confidence interval of predicted consumption expenditure of the country for the year 2125 would be $583.5116 billion and $686.4884 billion.

# EXERCISE

1. In a simple linear regression model, $Y_i = \alpha + \beta X_i + u_i$ for $i = ?$, why do we insert the random disturbance term $u_i$?

2. State and explain the assumptions of a classical linear regression model (CLRM)

3. In a simple linear regression model of the form $Y = \alpha + \beta X + u$, $i = ?$, n how can you estimate the regression parameters $\alpha$ and $\beta$?

4. Describe briefly the method of moments, used in estimating the regression parameters in a two variable linear regression model

5. Describe briefly the method of least squares used in estimating the regression parameters relating to a two variable linear regression model

6. How can you estimate a linear function (two variable) whose intercept is zero?

7. How can you estimate the elasticities from an estimated regression line?

8. State and prove the properties of the least squares estimators relating to a two variable linear regression model (CLRM)

9. Show that in a classical linear regression model the estimated regression parameters are unbiased.

10. Determine the mean and variance of $\alpha$ and $\beta$ relating to a model $Y = \alpha + \beta X + u$ for $i = 1, 2, \ldots, n$.

estimate ___ and $\beta$ and ___ calculate estimates ___ as errors of ___ estimates ___ to make the ___ conditional mean value of ___ corresponding to a ___ given ___ fixed of $X$ ___

**27.** The following table shows the ___ estimated expenditures and the ___ ___ interest rate ___ for the ten year period

| Year | ___ | ___ | ___ | ___ | ___ | ___ | ___ | ___ | ___ | ___ |
|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Investment | ___ | ___ | ___ | ___ | | | | | ___ | |
| Interest | ___ | ___ | ___ | ___ | | | | | | |

Test the hypothesis that investment is interest ___ by fitting a regression line to the above data and ___ indicating the ___ two levels of significance.

**28.** Given the following data

$\sum X_1 = ___$, $\sum X_2^2 = ___$, $\sum x_2 = 500$, $\sum ___ = ___$, $n = ___$ estimate the parameters

in the model $Y_1 = a + \beta X_1 + u_1$, and test the hypothesis $H_0: \beta = ___$ against the alternative (i) $\beta = 1.5$, (ii) $\beta > 1.5$, (iii) $\beta = 1.5$.

**29.** The true relationship between $X$ and $Y$ in the population is given by
$Y = 7 + 3X + u$. Suppose the value of $X$ in the sample of 10 observations are ___ , ___ , ___ The values of the disturbances are drawn at random from a normal population with zero mean and constant variance.

$u_1 = 0.464$, $u_2 = 0.06$, $u_3 = -1.48$, $u_4 = ___$, $u_5 = ___$, $u_6 = ___$, $u_7 = ___$, $u_8 = 0.64$, $u_9 = 0.18$ and $u_{10} = -1.37$.

(i) Present the 10 observed values of $X$ and $Y$

(ii) Estimate the least squares estimates of the regression coefficients and their standard errors.

(iii) Obtain the predicted value of $Y$ for $X = 12$.

**30.** The following data gives the production of coal and the number of wage earners in the coal industry

Output　　2.0.8　2.0　211.5　208.9　207.4　203.3　198.8　192.1　183.2　76.8　million tonnes.

Number of
workers　706.2　703　701.8　699.1　697.4　795.3　692.7　630.2　612.1　53.___ 000's)

(i) Estimate the production function (linear) of coal
(ii) Find average and marginal productivity of labour.
(iii) Estimate $t$-ratios and test their significance.

**31.** The following are data on
$Y = $ Quit rate per 100 employees in manufacturing
$X = $ unemployment rate
The data are for the United States and cover the period 1960-1972

| Year | Y | X | Year | Y | X |
|------|-----|-----|------|-----|-----|
| 1960 | 1.3 | 6.2 | 1966 | 2.6 | 3.7 |
| 1961 | 1.2 | 7.8 | 1967 | 2.3 | 3.6 |
| 1962 | 1.4 | 5.8 | 1968 | 2.5 | 3.3 |
| 1963 | 1.4 | 5.7 | 1969 | 2.7 | 3.3 |
| 1964 | 1.5 | 5.0 | 1970 | 2.1 | 5.6 |
| 1965 | 1.9 | 4.0 | 1971 | 1.8 | 6.8 |
| | | | 1972 | 2.2 | 5.6 |

11. ...

12. ...

13. ...

14. ...

15. ...

16. Show that the least squares estimators $\alpha$ and $\beta$ are such that

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{X}, \qquad \hat{\beta} = \frac{\sum x_i y_i}{\sum x_i^2}$$

17. Describe the testing procedure of the significance of the regression coefficients in the model $Y = \alpha + \beta X + u_i$, for $i = 1 \ldots n$

18. What is meant by goodness of fit of the correlation coefficient $R^2$?

19. Show that Total sum of squares = Explained sum of squares + unexplained sum of squares.

20. What is coefficient of determination? Show that it lies between 0 and 1 and hence show that the value of correlation coefficient between two variables lies between ... and ...

21. How can you formally write the regression results of the regression model $Y = \alpha + \beta X + u_i$, where $u_i$ (for $i = 1 \ldots n$) satisfies all the properties of CLRM?

22. How can you use the analysis of variance in the simple classical linear regression model.

23. What is the meaning of the term 'Prediction'? How can you incorporate the term in the CLRM? Distinguish between point prediction and interval prediction in this regard.

24. Show that the OLS point predictor in the CLRM satisfies the BLUE property.

25. The following sums were obtained from 16 pairs of observations on X and Y: $\Sigma X = 126$, $\Sigma X_i = 657$, $\Sigma X_i Y_i = 492$, $\Sigma Y_i = 63$, $\Sigma X_i = 96$. Estimate the parameters in the model $Y = \alpha + \beta X_i + u_i$ and $R^2$.

Test the hypothesis that $\beta = 2.0$

26. A sample of 20 observations corresponding to the regression model $Y = \alpha + \beta X + u$ gave the following data:

$\Sigma Y_i = 21.9$, $\Sigma_i Y_i$, $\bar{Y}^2 = 86.9$, $\Sigma(X - \bar{X})(Y_i - \bar{Y}_i) = 106.4$, $\Sigma X = 186.2$, $\Sigma_i X_i$ ...

32. ...

| Year | Income ... £m | Consumption ... |
|---|---|---|
| 196_ | ... | 4 _ |
| 196_ | _ _ | 8 |
| 196_ | ... | ... |
| 196_ | _ _ | ... |
| 200_ | _ 450 | 7 _ |

33. A random sample of ten families had the following income and food expenditure per week:

| Families | A | B | C | D | E | ... |
|---|---|---|---|---|---|---|
| Family income | | | | | | |
| Family expenditure | | | | | | |

Estimate the regression line of food expenditure on income and interpret your results.

34. The following results have been obtained from a sample of ... observations on the value of sales ( ) as a firm and the corresponding prices ( ).

$$\bar{X} = \ldots \quad \bar{Y} = \ldots \quad \sum XY = \ldots \quad \sum X = \ldots \quad \sum X^2 = \ldots$$

i. Estimate the regression line of sales on price and interpret the results.
ii. What is the part of the variation in sales which is not explained by the regression line?
iii. Estimate the price elasticity of sales.

35. The following table gives the quantities of commodity Z brought in each year ... and the corresponding prices.

| Year | 2009 | 2010 | 2011 | ... | ... | ... | ... | ... | ... | ... |
|---|---|---|---|---|---|---|---|---|---|---|
| Quantity (in tons) | 70 | 765 | 700 | 795 | 800 | 805 | ... | ... | 850 | ... |
| Price in £ | 13 | ... | 15 | ... | 7 | ... | | 7 | | |

i. Estimate the linear demand function for commodity Z.
ii. Calculate the price elasticity of demand.
iii. Forecast the demand at the mean price of the sample.
iv. Forecast the demand at P = 20.

36. A sample of 20 observations on a time series data on X and ... to be used for estimating the linear function $Y = a + bX + u$. The first 10 observations yield the following results:

$$\bar{X} = 15.30, \quad \bar{Y} = 160.00, \quad \sum_{i=1}^{10} \ldots \quad \sum y = \ldots 00, \quad \sum_{i=1}^{10} Y - \bar{Y}^2 = 45\,600$$

$$\sum_{i=1}^{10} X_i \ldots \bar{X}Y \bar{Y} = \ldots 598.00$$

... 10 ................ pairs of values of $X$ and $Y$ ..............

$\sum$ ............

... the regression ........... the .........

37. The following table ... data on the quantity supplied of a ... commodity
... and its price.

| Year | ... | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Quantity cost unit ... | ... | ... | ... | ... | ... | ... | ... | 40 | 30 |
| Export price $ ... | 10 | ... | ... | ... | ... | ... | ... | ... | ... |
| $ per unit) | | | | | | | | | |

(i) Test the hypothesis that the quantities supplied and price are related ... ... by estimating the export supply function ... ... ... ...

(ii) Show that $\beta$ is a part of the price elasticity of supply ... ... a numerical ... ... for the latter.

(iii) If price in year .... becomes $M$ and in .... $S$ ... then ... estimate the export volume of the commodity in these years.

38. The following table includes the total cost and the ... of ... of firm $X$ over a ten year period.

| Year | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 |
|---|---|---|---|---|---|---|---|---|---|---|
| Quantity $(X)$ (000 units) | 40 | 42 | 48 | 55 | 60 | 70 | 88 | 100 | 120 | 140 |
| Total cost $(Y)$ (000 dollars) | 150 | 140 | 160 | 170 | 150 | 162 | 185 | 165 | 190 | 155 |

(i) Estimate the linear cost function $Y = \alpha + \beta X$

(ii) Find the AVC, MC and AC and plot them roughly on a graph.

39. The total investment function for the economy as a whole is assumed to be of the form

$$I = \alpha r^\beta e^u$$

where $I$ = investment, $r$ = rate of interest

The following sample is given

| $I$ ($ billion) | 9.0 | 5.5 | 8.5 | 4.0 | 3.5 | 2.5 | 3.0 | 1.3 | 1.2 | 1.8 | 1.5 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $r$ (percent) | 2 | 3 | 2 | 4 | 4 | 6 | 4 | 6 | 6 | 7 | 9 |

(i) Estimate the parameters of the investment function by OLS

(ii) Test the statistical significance of the coefficients at 5% level of significance.

(iii) Construct a 95% confidence interval for $\beta$.

(iv) Find the value of $R^2$ and interpret the result.

40. a. Write down the assumptions essential for each of the following tasks

i Proving that the OLS estimators are unbiased

... ... ...

**41.** Consider the following estimated regression equation $\hat{Y} = \ldots + \ldots \hat{\beta} \ldots$ with their standard error of $\hat{\beta} \ldots$ It has further been given that $\hat{\beta} \ldots$

$$\sum_{i=1}^{n} \ldots \ldots$$

Find out the following:

(a) Sample size (b) The estimated intercept (c) Total sum of squares

(d) Residual sum of squares (RSS) (e) Estimated error variance $\hat{\sigma}_u$

**42.** (a) In a two variable regression model, show that $R^2 = \dfrac{SS}{TSS}$

(b) Explain why random error term is introduced in an econometric model.

(c) Which one would you consider to be CLRM ? First regression

(i) $Y_i = \beta_1 + \beta_2 X_i^2$   (ii) $Y_i = \beta_1 + \sqrt{\beta_2} X_i$   (iii) $Y_i = \beta_1 + \beta X_i^\beta$   (iv) $Y_i = \beta_1 + \beta_2 X_i$

**43.** Consider the following estimated two variable LRM $\hat{Y}_i = a_1 + 0.5 X_i$, with $n$

$$\bar{Y} = 0,\ \bar{Y} = 4,\ \sum Y_i^2 = 2201,\ \sum Y_i^2 = 4951$$

(a) Obtain the estimated regression coefficient when $X$ is regressed on $Y$

(b) Obtain the coefficient of correlation between $X$ and $Y$

(c) Obtain the unbiased estimate of the error variance when $Y$ is regressed on $X$

(d) Obtain the estimated value of the intercept term and its estimated standard error when $Y$ is regressed on $X$

(e) Test the suggestion that $Y$ is positively related to $X$ at 5% level of significance.

**44.** Consider the following regression equation $Y_i = \alpha + \beta X_i + u_i$, where $n = \ldots$

$$\sum Y_i = 80,\ \sum X_i^2 = 600,\ \sum Y_i^2 = 734,\ \sum X_i Y_i = 480$$

(a) Obtain the estimated value of $\alpha$ and $\beta$.

(b) Test the hypothesis that $X$ and $Y$ are negatively correlated against the hypothesis that they are not at 5% level of significance.

# 3
# Multiple Linear Regression Model

## 1.1 Introduction

In simple regression analysis we study the relationship between an explained dependent variable $Y$ and an explanatory independent variable $X$. In multiple regression analysis we study the relationship between $Y$ and a number of explanatory variables $X_1, X_2, \ldots, X_k$. For example, in demand studies, we may be interested in investigating the relationship between quantity demanded of a good, and prices of that good, prices of substitute goods and income of the consumer. In fact this problem can be analysed with the help of multiple regression analysis.

Let us consider a linear regression model where there are $k$ independent variables say $X_1, X_2, \ldots, X_k$ and $Y$ is the only dependent variable. In this case the regression model is given by,

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3i} + \cdots + \beta_K X_{Ki} + u_i, \text{ where } \qquad i \qquad n \qquad 3.1)$$

Here $\beta_0, \beta_1, \beta_2, \ldots, \beta_K$ are $(K + 1)$ regression parameters
[$K$ = number of explanatory variables and $K + 1$ = number of regression parameters
$u$ = random disturbance term = error term.
$\beta_0$ = constant term and $\beta_1, \beta_2, \ldots, \beta_K$ are the partial regression coefficients

**We make the following assumptions about $u_i$:**

(i) $E(u_i) = 0$ for all $i$, $i = 1, 2, \ldots, n$

(ii) $\text{var}(u_i) = \sigma_u^2$ for all $i$

(iii) $u_i$ and $u_j$ are independent for all $i \neq j$

(iv) $u_i$ and $Y$ are independent for all $i$ and $j$

(v) $u_i$ is normally distributed for all $i$ $[u_i \sim N(0, \sigma_u^2)]$

(vi) There is no linear dependencies in the explanatory variables.

Under the first four assumptions, we can show that the method of least squares gives estimators of $\beta_0, \beta_1, \beta_2, \ldots, \beta_K$ that are unbiased and have minimum variance.

In equation (3.1) if we put $i = 1, 2, \ldots, n$, we have

For $i = 1$, $Y_1 = \beta_0 + \beta_1 X_1 + \beta_2 X_{21} + \cdots \beta_K X_{K1} + u_1$

$i = 2$, $Y_2 = \beta_0 + \beta_1 X_{12} + \beta_2 X_{22} + \cdots + \beta_K X_{K2} + u_2$

$i = 3$, $Y_3 = \beta_0 + \beta_1 X_{13} + \beta_2 X_{23} + \cdots \beta_K X_{K3} + u_3$

$i = n$, $Y_n = \beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \cdots + \beta_K X_{Kn} + u_n$

**105**

in which the classical set of $k$ equations can be written as

$$Y + X\beta + \varepsilon \qquad (3.2)$$

where

Equation represents a set of $n$ equations where there are $k$ independent explanatory variables with a finite of sample size.

The model is said to be a Classical Linear Regression Model (CLRM) if it satisfies the following properties

i. $E(u)$ is a null vector

where $u = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix}$ $E(u) = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

since $E(u_i) = 0$ for all $i = 1 \cdots n$

ii. The dispersion matrix, variance-covariance matrix of a disturbance vector is a scalar matrix

$D(u) = \sigma_u^2 I_n$ where $I_n = \begin{vmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & & & \\ 0 & 0 & \cdots & 1 \end{vmatrix}$ = identity matrix and $\sigma_u^2$ is a constant

**Proof** Let us suppose $u$ is a random vector variable then dispersion matrix, $D(u) = E[u - E(u)][u - E(u)]'$

$= E(uu')$ since $E(u) = 0$

$E \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{pmatrix} \begin{pmatrix} u_1 & u_2 & \cdots & u_n \end{pmatrix}$

$= E \begin{bmatrix} u_1^2 & u_1 u_2 & \cdots & u_1 u_n \\ u_2 u_1 & u_2^2 & \cdots & u_2 u_n \\ u_3 u_1 & u_3 u_2 & \cdots & u_3 u_n \\ \vdots & & & \\ u_n u_1 & u_n u_2 & \cdots & u_n^2 \end{bmatrix} = \begin{bmatrix} E(u_1^2) & E(u_1 u_2) & \cdots & E(u_1 u_n) \\ E(u_2 u_1) & E(u_2^2) & \cdots & E(u_2 u_n) \\ E(u_3 u_1) & E(u_3 u_2) & \cdots & E(u_3^2) \\ \vdots & & & \\ E(u_n u_1) & E(u_n u_2) & \cdots & E(u_n^2) \end{bmatrix}$

$$
\begin{bmatrix}
\sigma_u^2 & 0 & & & 0 \\
0 & \sigma_u^2 & \cdots & \cdots & 0 \\
\vdots & & \ddots & & \vdots \\
0 & & \cdots & & \sigma_u^2
\end{bmatrix}
$$

Since $E(u_i u_j) = 0$ for $i \neq j$, $i \neq j$

and $E \cdots \cdots = \sigma^2$ for all ...

$$
\text{Var}(u) = \sigma_u^2 
\begin{bmatrix}
1 & 0 & & & 0 \\
0 & 1 & & & 0 \\
\vdots & & \ddots & & \vdots \\
0 & & & & 1
\end{bmatrix} = \sigma_u^2 I_n.
$$

$u \sim N_n I_n$

$u$ is a multivariate normal vector with mean $0_{n \times 1}$ and variance co-variance matrix $\sigma_u^2 I_n$

• $X$ – the independent variables $X_1, X_2 \ldots X_k$ are non-stochastic or nonrandom or $X$ is a non-stochastic matrix.

• Rank of matrix $X$ is $(K + 1)$

Rank of a matrix implies the maximum number of linearly independent columns of the matrix

Since $X$ is a matrix of order $n \times (K + 1)$, all the columns of $X$ should be linearly independent. Now $X'$ is a matrix of order $(K + 1) \times n$ and $(X'X)$ is a matrix of order $(K + 1)(K + 1)$. If the rank of $X$ is $(K + 1)$ then rank of $(X'X)$ is also $(K + 1)$, and $|X'X| \neq 0$. If its rank is $(K + 1)$ and $|X'X| = 0$ if its rank is $< (K + 1)$.

If $|X'X| = 0$ then $(X'X)^{-1}$ does not exist.

## 3.2. The Least Squares Method (OLS) for Estimation of Regression Parameters

In vector matrix form the general linear regression model (Equation 3.1) can be written as $Y = X\beta + u$

where $Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_{n \times 1} \end{bmatrix}$, $X = \begin{bmatrix} 1 & X_{11} & X_{21} & \cdots & X_{k1} \\ 1 & X_{12} & X_{22} & \cdots & X_{k2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{1n} & X_{2n} & \cdots & X_{kn} \end{bmatrix}_{n \times (k+1)}$

$u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}_{n \times 1}$ and $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix}_{(K+1) \times 1}$

Let $\hat{Y} = X\hat{\beta}$ be the vector of the regressed value of $Y$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}$$

and $\hat{\beta}$ be the vector of estimators

where $\hat{Y}$ is a $n \times 1$ order vector and $X$ is a $n \times (K+1)$ order matrix

Let $e$ be the residual vector i.e. $e = Y - \hat{Y}$ where $\hat{Y} = X\hat{\beta}$ and $Y = Y$

$$e = Y - X\hat{\beta}$$

Here, $e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$ and $e' = [e_1 \ e_2 \ \cdots \ e_n]$

Now $e'e = [e_1 \ e_2, \ , e_n] \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = e_1^2 + e_2^2 + \cdots + e_n^2 = \sum_{i=1}^{n} e_i^2$ i.e. $e'e = \sum_{i=1}^{n} e_i^2$

$$e'e = (Y - X\hat{\beta})'(Y - X\hat{\beta})$$

$$= Y'Y - \hat{\beta}'X'Y - Y'X\hat{\beta} + \hat{\beta}'X'X\hat{\beta}$$

Here, $\hat{\beta}'X'Y$ is scalar $(1 \times 1)$, it is equal to its transpose i.e. $\hat{\beta}'X'Y = Y'X\hat{\beta}$

Now by OLS method we have to minimise $\sum_{i=1}^{n} e_i^2 = e'e$ with respect to $\hat{\beta}$

Now $\dfrac{d(e'e)}{d\hat{\beta}} = 0 \cdot X'Y - X'Y + 2X'X\hat{\beta} = 0$

or $2X'X\hat{\beta} = 2X'Y$ or $X'X\hat{\beta} = X'Y$

or, $\hat{\beta} = (X'X)^{-1}X'Y$, where $\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_K \end{bmatrix}_{(K+1) \times 1}$

**To derive this result more clearly we consider a three variable (with two explanatory variables i.e., when $K = 1$) linear regression model which takes the form**

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \ i = 1, 2, \ldots, n$$

Let us define $X =$

in vector-matrix form the equation,

can be written as $Y = X\beta + u$

Now $X'X =$

Now $\beta'X'X\beta =$

$$\beta'X'X\beta = \beta_1^2 \Sigma x_{1i}^2 + 2\beta_1\beta_2 \Sigma x_{1i} x_{2i} + \beta_2^2 \Sigma x_{2i}^2$$

Now $\dfrac{d}{d\beta}(\beta'X'X\beta) = \begin{bmatrix} \dfrac{d}{d\beta_1}(\beta'X'X\beta) \\ \dfrac{d}{d\beta_2}(\beta'X'X\beta) \end{bmatrix}$

$$\frac{d}{d\beta}(Y'X\beta + \beta'X'Y)$$

Again $Y'X\beta = \begin{bmatrix} y_1 & \cdots & y_n \end{bmatrix} \begin{vmatrix} & & \\ & & \beta_1 \\ & & \beta \end{vmatrix}$

$$= \begin{vmatrix} \Sigma x_{1i} y_i & \Sigma x_{2i} y_i \end{vmatrix} \begin{vmatrix} \beta_1 \\ \beta_2 \end{vmatrix} = \beta_1 \Sigma x_{1i} y_i + \beta_2 \Sigma x_{2i} y_i$$

Now $\dfrac{d}{d\beta}(Y'X\beta) = \begin{bmatrix} \dfrac{d}{d\beta_1}(Y'X\beta) \\[2mm] \dfrac{d}{d\beta_2}(Y'X\beta) \end{bmatrix} = \begin{bmatrix} \Sigma x_{1i} y_i \\ \Sigma x_{2i} y_i \end{bmatrix} = X'Y$

$$\frac{d}{d\beta}(Y'X\beta) = X'Y$$

Again $\beta'X'Y = \begin{bmatrix} \beta_1 & \beta_2 \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$

$$= \begin{bmatrix} \beta_1 & \beta_2 \end{bmatrix} \begin{bmatrix} \Sigma x_{1i} y_i \\ \Sigma x_{2i} y_i \end{bmatrix} = \beta_1 \Sigma x_{1i} y_i + \beta_2 \Sigma x_{2i} y_i$$

Now $\dfrac{d}{d\beta}(\beta'X'Y) = \begin{bmatrix} \dfrac{d}{d\beta_1}(\beta'X'Y) \\[2mm] \dfrac{d}{d\beta_2}(\beta'X'Y) \end{bmatrix} = \begin{bmatrix} \Sigma x_{1i} y_i \\ \Sigma x_{2i} y_i \end{bmatrix} = X'Y$

Thus we have $\dfrac{d}{d\beta}(\beta'X'Y\beta) = 2X'X\beta$

$$\frac{d}{d\beta}(Y'X\beta) = X'Y \quad \text{and} \quad \frac{d}{d\beta}(\beta'X'Y) = X'Y$$

Since $e'e = Y'Y - \beta'X'Y - Y'X\beta + \beta'X'X\beta$

… we can … $\hat{\beta} = \dots$

$$\hat{\beta} = \begin{bmatrix} \Sigma x_{1i}^2 & \Sigma x_{1i} x_{2i} \\ \Sigma x_{1i} x_{2i} & \Sigma x_{2i}^2 \end{bmatrix} \begin{bmatrix} \Sigma x_{1i} y_i \\ \Sigma x_{2i} y_i \end{bmatrix}$$

For A … with two explanatory variables … by solving the equations we can find out the values of … and β

Now $\hat{\beta} = (X'X)^{-1} X'Y$

$$= \frac{1}{|X'X|} \; adj\,(X'X)\, X'Y$$

Now Adj $(X'X)$ = Transpose of cofactor matrix … $(X'X)$

$$= \begin{bmatrix} \Sigma x_2^2 & \Sigma x_{1i} x_{2i} \\ -\Sigma x_{1i} x_{2i} & \Sigma x_1^2 \end{bmatrix} \quad \text{and} \quad \dots$$

$$\text{Adj}\,(X'X) = \begin{bmatrix} \Sigma x_{2i}^2 & -\Sigma x_{1i} x_{2i} \\ -\Sigma x_{1i} x_{2i} & \Sigma x_{1i}^2 \end{bmatrix}$$

$$\begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = \frac{1}{|X'X|} \begin{bmatrix} \Sigma x_{2i}^2 & -\Sigma x_{1i} x_{2i} \\ -\Sigma x_{1i} x_{2i} & \Sigma x_{1i}^2 \end{bmatrix} \begin{bmatrix} \dots \\ \dots \end{bmatrix}$$

$$\hat{\beta}_1 = \frac{1}{|X'X|} (\Sigma x_2^2 \Sigma x_{1i} - \Sigma x_{1i} x_i \Sigma x_{2i} \dots) = \dots$$

and $\beta_2 = \frac{1}{|X'X|} (\Sigma x_{1i} x_{2i} \Sigma x_{1i} y - \Sigma x_{1i}^2 \Sigma \dots) = \dots$

When $\hat{\beta}_1$ and $\beta_2$ are known, $\beta_0$ can be obtained from the relation

$$\bar{Y} = \beta_0 + \beta_1 \bar{X} + \beta_2 \bar{X}_2, \quad \beta_0 = \bar{Y} - \beta_1 \bar{X}_1 - \beta_2 \bar{X}_2$$

**We can also find out the values of $\beta_1$ and $\beta_2$ directly by using Cramer's rule**

Since $\beta = (X'X)^{-1} X'Y$ or $(X'X)\beta = X'Y$

or 
$$\begin{vmatrix} \Sigma n_i & \Sigma ... & \mu \\ \Sigma n_i & \Sigma ... & \mu \end{vmatrix}$$

or,
$$\begin{vmatrix} \Sigma ... & \Sigma ... & \beta \\ \Sigma ... & \Sigma ... & \mu \end{vmatrix}$$

or   $\beta_1 \Sigma x_{1i}^2 + \beta_2 \Sigma x_{1i} x_{2i} = \Sigma x_{1i} y_i$ ——— (A

and  $\beta_1 \Sigma x_{1i} x_{2i} + \beta_2 \Sigma x_{2i}^2 = \Sigma x_{2i} y_i$ ——— (B)

Solving equations A and B by Cramer's rule
we have

$$\beta_1 = \frac{\begin{vmatrix} \Sigma x_{1i} y & \Sigma x_{1i} x_{2i} \\ \Sigma x_{2i} y & \Sigma x_{2i}^2 \end{vmatrix}}{\begin{vmatrix} \Sigma x_{1i}^2 & \Sigma x_{1i} x_{2i} \\ \Sigma x_{1i} x_{2i} & \Sigma x_{2i}^2 \end{vmatrix}} = \frac{\Sigma x_{1i} y \, \Sigma x_{2i}^2 - \Sigma x_{2i} y \, \Sigma x_{1i} x_{2i}}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$$

and  
$$\beta_2 = \frac{\begin{vmatrix} \Sigma x_{1i}^2 & \Sigma x_{1i} y \\ \Sigma x_{1i} x_{2i} & \Sigma x_{2i} y \end{vmatrix}}{\begin{vmatrix} \Sigma x_{1i}^2 & \Sigma x_{1i} x_{2i} \\ \Sigma x_{1i} x_{2i} & \Sigma x_{2i}^2 \end{vmatrix}} = \frac{\Sigma x_{1i}^2 \, \Sigma x_{2i} y - \Sigma x_{1i} x_{2i} \, \Sigma x_{1i} y}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$$

When $\beta_1$ and $\beta_2$ are known $\beta_0$ is obtained from the relation

$$\bar{Y} = \beta_0 + \beta_1 \bar{X}_1 + \beta_2 \bar{X}_2 \qquad \beta_0 = \bar{Y} - \beta_1 \bar{X}_1 - \beta_2 \bar{X}_2$$

**Note**   For the three variable linear regression equation

$(Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad i = 1, 2, \ldots n$ where $u_i \sim N(0, \sigma_u^2))$
we can also find out the values of the regression parameters in another way
This method is described below.
The estimated regression line is given by

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} \quad \text{and} \quad \bar{Y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X}_1 + \hat{\beta}_2 \bar{X}_2$$

where $\hat{\beta}_0$, $\hat{\beta}_1$ and $\hat{\beta}_2$ are the OLS estimators of $\beta_0$, $\beta_1$ and $\beta_2$

We obtain by subtraction $y_i = Y_i - \bar{Y}$

$$= \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} - \hat{\beta}_0 - \hat{\beta}_1 \bar{X}_1 - \hat{\beta}_2 \bar{X}_2$$

$\beta_1 x_{1i} + \beta_2 x_{2i} + \bar{Y}$

$\beta_{0} \ldots \beta \ldots$ where $a = 1$ $\bar{U}$ $= 1$ $\bar{Y}$

Now errors of estimate $e_i = y_i \ldots = \epsilon - (\beta x \ldots \beta)$

and $\Sigma e_i^2 = \Sigma (y_i - \beta_1 x_{1i} - \beta_2 x_{2i})^2$

The first order conditions for minimization require

$$\frac{\partial \Sigma e^2}{\partial \beta_1} \quad 2 \Sigma y_i (\beta_1 x_{1i} \quad \beta_2 x_{2i}, \quad a_{1i}) = 0$$

i.e. $\Sigma x_{1i} y_i = \Sigma x_{1i}^2 + \beta \Sigma x_{1i} x_{2i}$ (1)

$$\frac{\partial \Sigma e^2}{\partial \beta_2} \quad 2 \Sigma y_i \quad \hat{\beta} x_{1i} \quad \beta_2 x \quad \ldots e_2 \quad 0$$

iii. $\Sigma x_2 y_i = \beta_1 \Sigma x_{1i} x_{2i} + \beta_2 \Sigma x_2^2$ (2)

Now solving equations (1) and (2) by Cramer's rule we have,

$$\beta_1 = \div \frac{\begin{vmatrix} \Sigma y_i & \Sigma x_2 \\ \Sigma x_2 y & \Sigma x_2^2 \end{vmatrix}}{\begin{vmatrix} \Sigma x_1^2 & \Sigma x_2 \\ \Sigma x_1 x_2 & \Sigma x_2^2 \end{vmatrix}} \quad \frac{\Sigma x_2^2 \Sigma x_{1i} y_i \quad \Sigma x_{1i} x_2 \quad \Sigma x_2 y}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2}$$

and $\beta_2$

$$\frac{\begin{vmatrix} \Sigma x_1^2 & \Sigma x_1 y \\ \Sigma x_1 x_2 & \Sigma x_2 y \end{vmatrix}}{\begin{vmatrix} \Sigma x_1^2 & \Sigma x_1 x_2 \\ \Sigma x_1 x_2 & \Sigma x_2^2 \end{vmatrix}} + \frac{\Sigma x_1^2 \Sigma x_2 y \quad \Sigma x_1 x_2 \quad \Sigma x_1 y}{\Sigma x_1^2 \Sigma x_2^2 \quad (\Sigma x_1 x_2)^2}$$

When $\beta_1$ and $\beta_2$ are known, $\beta_0$ can be obtained from the relation

$$\bar{Y} = \beta_0 + \beta_1 \bar{X}_1 + \beta_2 \bar{X}_2 \quad \text{i.e.} \quad \beta_0 \quad \bar{Y} \quad \beta_1 \bar{X}_1 \quad \beta_2 \bar{X}_2$$

**Example 3.1** Consider the following regression model $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u_i$

where $u_i$ is normally distributed with mean 0 and variance $\sigma_u^2$

| Y | 4 | 7 | 3 | 9 | 7 |
|---|---|---|---|---|---|
| X | 2 | 3 | 1 | 5 | 9 |
| $X_2$ | 5 | 3 | 2 | 1 | 7 |

Estimate $\beta_0$, $\beta_1$ and $\beta_2$ (the OLS estimators of $\beta_0$, $\beta_1$ and $\beta_2$)

62-E

**Solution** Calculations for the estimators of the regression parameters

Table 3.1

Here n = 5 as five sets of values of the variables are given.

We know that the values of $\beta_0$ and $\beta$ in the regression equation

$Y = \beta_0 + \beta_1 X + \beta_2 X + u$ can be obtained by the O.L.S. method. The estimators of $\beta$ and $\beta$ are $\beta_1$ and $\beta_2$.

where $\beta_1 = \dfrac{\Sigma x_1 \Sigma x_2 y - \Sigma x_2 \Sigma x_1 y}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2}$

$= \dfrac{464 - 5}{926 \cdot 269} = \dfrac{5.4}{619} = 0.0059 \approx 0.806$

and $\beta_2 = \dfrac{\Sigma x_1^2 \Sigma x_2 y - \Sigma x_1 x_2 \Sigma x_1 y}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2} = \dfrac{40 + 3.1 + 20}{40 + 23.20 \cdot (17)^2}$

$= \dfrac{21 \cdot 340}{926 \cdot 269} = \dfrac{464}{619} = 0.198 = 0.720$

$\beta_1 = 0.806$ and $\beta_2 = 0.720$

When $\beta_1$ and $\beta_2$ are known $\beta_0$ can be obtained from the relation

$\bar{Y} = \beta_0 + \beta_1 X + \beta_2 X$

$\beta_0 = \bar{Y} - \beta_1 \bar{X} - \beta_2 \bar{X} = 6 - 0.806 \times 4 - 0.720 + 1.6 = 5.368$

$\beta_0 = 5.368,\ \beta_1 = 0.806$ and $\beta_2 = 0.720$

### 3.2.1 The Regression Coefficients Expressed in terms of Variances (SDs) and Coefficient of Correlations

In a three variable linear regression model $Y = \beta_0 + \beta_1 X + \beta_2 X + u$, $u_{r,r}$
$u_r$

Now putting values in the expression of $\beta_1$ we get

$$\beta_1 = \frac{\sigma_Y}{\sigma_X}\left[ \quad \right]$$

Proceeding in the same way we get $\Sigma \beta^{\ldots}$

$$\beta \ldots \qquad \text{where the symbols have the usual meaning}$$

### 3.2.2 Determination of Variances and Covariances of the Estimators of the Regression Parameters in Three Variable Linear Regression Model

For a three variable linear regression model, we assume the regression equation of the form

$$= \beta_0 + \beta_1 X_{1i} + \beta_2 X_2 \ldots \qquad i = 1 \ldots n$$

$$+ \beta_0, \beta_1 X \quad \beta_2 \bar{X} \quad + u \quad \text{where } u = 0$$

Now $\quad Y, \bar{Y} \quad y_i \quad Y \quad \bar{Y} \quad = \beta \quad X \quad \bar{X} + u_i \ldots$

or $\quad = \beta_1 x_{1i} \quad \beta_2 x_{2i} \quad u_i \quad \text{where } x_{ji} = X \quad \ldots \quad \bar{X} \quad \text{and } u = $

In vector-matrix form, the set of $n$ equations for $i \ldots n$ can be written in

$$Y = X\beta + u$$

where $X$ ...... and $\beta$ ......

Now $\quad X'X = \begin{bmatrix} \ldots \end{bmatrix}$

Now $\quad E(\hat\beta = \begin{bmatrix} E \hat\beta_1 & \ldots \\ E(\hat\beta_1 & \ldots \end{bmatrix}$

since $E \hat\beta = \beta$ and $\hat\beta = \beta$

$$\begin{bmatrix} \text{var}(\hat\beta_1) & \text{cov}(\hat\beta_1,\hat\beta_2) \\ \text{cov}(\hat\beta_1,\hat\beta_2) & \text{var}(\hat\beta_2) \end{bmatrix} = \sigma_u^2 (X'X)^{-1}$$

[See Property 2 of O.L.S estimator vector]

$$= \sigma_u^2 \begin{bmatrix} \Sigma x_1^2 & \Sigma x_1 x_2 \\ \Sigma x_1 x_2 & \Sigma x_2^2 \end{bmatrix}$$

$$= \sigma_u^2 \frac{\text{Adj} \begin{bmatrix} \Sigma x_1^2 & \Sigma x_1 x_2 \\ \Sigma x_1 x_2 & \Sigma x_2^2 \end{bmatrix}}{\begin{vmatrix} \Sigma x_1^2 & \Sigma x_1 x_2 \\ \Sigma x_1 x_2 & \Sigma x_2^2 \end{vmatrix}}$$

$$= \begin{bmatrix} n_u & \Sigma x_1 & \Sigma x_2 \\ \Sigma x_1 & \Sigma x_1 x_2 & \Sigma x_1 x_2 & \Sigma x_2^2 \\ \Sigma x_1 x_2 & \Sigma x_2^2 \end{bmatrix}$$

Now,
$$\begin{bmatrix} var(\beta_1) & cov(\beta_1,\beta_2) \\ cov(\beta_1,\beta_2) & var(\beta_2) \end{bmatrix} = \Sigma x_2^2 \cdot (\cdots) = \begin{bmatrix} \sigma_u^2 & \Sigma x & \Sigma x \\ \Sigma x_1 & \Sigma x \end{bmatrix}$$

$$= \sigma_u \begin{bmatrix} \dfrac{\Sigma x_2^2}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2} & \dfrac{\Sigma x_1 x_2}{\cdots} \\ \dfrac{\Sigma x_1 x_2}{\cdots} & \dfrac{\Sigma x_1^2}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2} \end{bmatrix}$$

$$var(\beta_1) = \frac{\sigma_u^2 \Sigma x_2^2}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2} = \frac{\sigma_u^2 n \sigma_{x_2}^2}{n \sigma_{x_1}^2 \sigma_{x_2}^2 \cdots}$$

$$var(\beta_1) = \frac{\sigma_u^2 n \sigma_{x_2}^2}{n \sigma_{x_1}^2 n \sigma_{x_2}^2 \cdots} = \frac{\sigma_u^2}{n \sigma_{x_1}^2 (1 - r_{12}^2)}$$

[Since $\Sigma x_1^2 = \Sigma (X_1 - \bar{X}_1)^2 = n \frac{1}{n} \Sigma (X_1)^2$ ... $n \sigma_{x_1}^2$ ... similarly, $\Sigma x_2^2 = n \sigma_{x_2}^2$ ...

and $\Sigma x_1 x_2 = \Sigma X_1 \bar{X} \cdots = n \frac{1}{n} \Sigma (\cdots) = n \, cov(X_1, X_2) = n \, r_{12} \sigma_{x_1} \sigma_{x_2}$ as $\frac{cov(X_1, X_2)}{\sigma_{x_1} \sigma_{x_2}} = r_{12}$ ]

Similarly, $var(\beta_2) = \dfrac{\sigma_u^2 \Sigma x_1^2}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2}$

$$= \frac{\sigma_u^2 n \sigma_{x_1}^2}{n \sigma_{x_1}^2 n \sigma_{x_2}^2 - n^2 r_{12}^2 \sigma_{x_1}^2 \sigma_{x_2}^2} = \frac{\sigma_u^2 n \sigma_{x_1}^2}{n \sigma_{x_1}^2 n \sigma_{x_2}^2 (1 - r_{12}^2)}$$

$$var(\beta_2) = \frac{\sigma_u^2}{n \sigma_{x_2}^2 (1 - r_{12}^2)}$$

Again, $cov(\beta_1, \beta_2) = \dfrac{-\sigma_u^2 \Sigma x_1 x_2}{\Sigma x_1^2 \Sigma x_2^2 - (\Sigma x_1 x_2)^2}$

$$= \frac{-\sigma_u^2 n r_{12} \sigma_{x_1} \sigma_{x_2}}{n \sigma_{x_1}^2 n \sigma_{x_2}^2 (1 - r_{12}^2) \sigma_{x_1} \sigma_{x_2}} = \frac{-n \sigma_u^2 r_{12}}{n \sigma_{x_1} \sigma_{x_2} (\cdots)}$$

Now, $\mathrm{var}(\beta_0 + \beta_1 \dots) = \mathrm{var}(\beta_0) + \dots \mathrm{var}(\beta_1) + 2\,\mathrm{cov}(\beta_0, \beta_1)$

$$\frac{\sigma_u}{\dots} \qquad \sigma_u \qquad \dots$$

and $\mathrm{var}\,\beta_1 \quad (\beta_0) + \mathrm{var}(\beta_0) + \mathrm{var}(\beta_0) + 2\,\mathrm{cov}(\beta_0)$

$$-\frac{\sigma_u}{n\sigma_x} \qquad \sigma_u \qquad \dots$$

since $\bar{Y} = \beta_0 + \beta_1 \bar{X} + \beta \bar{X} \qquad \beta_0 = \dots \quad \beta_1 \dots$

It can be seen that

$$\mathrm{var}(\beta_0) = \frac{\sigma_u^2}{n} + \bar{X}\,\mathrm{var}(\beta) + 2\bar{X}\,\bar{X}\,\mathrm{cov}(\beta_1,\beta_2) \quad \bar{X}\,\mathrm{var}(\beta_2)$$

$$\mathrm{cov}(\beta_0,\beta_2) = -\bar{X}\,\mathrm{var}(\beta) + \bar{X}\,\mathrm{cov}(\beta_1,\beta_2)$$

and $\mathrm{cov}(\beta_0,\beta_1) = \{\bar{X}\,\mathrm{cov}(\beta_1,\beta_2)\} \quad \bar{X}\,\mathrm{var}(\beta_1)\}$

**Note** In calculating $\mathrm{var}(\beta_0)$, $\mathrm{var}(\beta_1)$, $\mathrm{var}(\beta_2)$, $\mathrm{cov}(\beta_1,\beta_2)$, $\mathrm{cov}(\beta_0,\beta_1)$ and $\mathrm{cov}(\beta_0,\beta_2)$ if $\sigma_u^2$ is not known it is to be replaced by its unbiased estimator $\sigma_u = \Sigma e^2 \quad n$

**Example 3.2.** The following table presents data on a sample of five persons randomly drawn from a large firm giving their annual salaries in thousands of dollars $Y$, years of education $X$, and years of experience with the firm they are working $(X_2)$

| $Y$ | 30 | 30 | 28 | 34 | 40 |
|---|---|---|---|---|---|
| $X$ | 4 | 7 | 6 | 4 | 5 |
| $X_2$ | 10 | 8 | 11 | 9 | 17 |

Assuming a linear regression of the form

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad u_i \sim N(0, \sigma_u^2)$$

(i) find the OLS estimators $\beta_0$, $\beta_1$ and $\beta_2$

(ii) find the value of $r_{XY}$

(iii) find the estimated regression equation.

(iv) find $\Sigma e_i^2$

(v) find the values of $\mathrm{var}(\beta_1)$, $\mathrm{var}(\beta_2)$ and $\mathrm{var}(\beta_0)$

(vi) find the value of $\mathrm{cov}(\beta_1,\beta_2)$

**Solution :**

**Calculation Table 3.2**

| Y | $X_{1i}$ | $X_{2i}$ | $y_i$ $= Y_i - \bar{Y}$ | $y_i^2$ | $x_{1i}$ $= Y_{1i} - \bar{X}_1$ | $x_{1i}^2$ | $x_{2i}$ $= X_{2i} - \bar{X}_2$ | $x_{2i}^2$ | $y_i x_{2i}$ | $y x_{1i}$ | | | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 30 | 4 | 10 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 6 | | | |
| 20 | 3 | 8 | -10 | 100 | 2 | 4 | 2 | 4 | 20 | 20 | 4 | | |
| 36 | 6 | 11 | 6 | 36 | 1 | 1 | 1 | 1 | 6 | 6 | | | |
| 24 | 4 | 9 | -6 | 36 | 1 | 1 | 1 | 1 | 6 | 6 | | | |
| 40 | 8 | 12 | 10 | 100 | 3 | 9 | 2 | 4 | 40 | 30 | | | |

| $\Sigma Y_i$ | $\Sigma X_{1i}$ | $\Sigma X_{2i}$ | $\Sigma y_i$ | $\Sigma y_i^2$ | $\Sigma x_{1i}$ | $\Sigma x_{1i}^2$ | $\Sigma x_{2i}$ | $\Sigma x_{2i}^2$ | $\Sigma x$ | $\Sigma$ | $\Sigma$ | |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| = 50 | 25 | = 50 | = 0 | = 272 | = 0 | = 16 | = 0 | = 10 | 6 | | 2 | |

Here $n = 5$ as five sets of values are given.

Now $\bar{Y} = \dfrac{\Sigma Y_i}{n} = \dfrac{150}{5} = 30, \quad \bar{X}_1 = \dfrac{\Sigma X_{1i}}{n} = \dfrac{25}{5} = 5, \quad \bar{X}_2 = \dfrac{\Sigma X_{2i}}{n} = \dfrac{50}{5} = 10$

We have to find out the OLS estimators $\beta_0$, $\beta_1$ and ...

We know that $\beta_1$ ...

We now put the values from the calculation table and get

$$\beta_1 = \ldots = \frac{6.50 - 6.24}{60 - 14.4} = \frac{4}{6} = -0.25$$

Similarly $\beta_2$ ...

When $\beta_1$ and $\beta_2$ are known, $\beta_0$ can be obtained from the relation

$$\beta_0 = \bar{Y} - \beta_1 \ldots$$

$$\beta_0 = 30 - \ldots$$

Thus the OLS estimators of the parameters are $\beta_0 = 23.75$, $\beta_1 = \ldots$ and ...

(ii) We have to find out the value of product moment correlation coefficient between the two explanatory variables $X_1$ and $X_2$ i.e. $r_{12}$

We know that $r_{12} = \ldots$

$$= \ldots = 0.5 \quad r_{12} = 0.5$$

(iii) The estimated regression line is given by

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i}$$

$$\hat{Y}_i = 23.75 - 0.25 X_{1i} + 5.5 X_{2i} \text{ is the estimated regression line equation}$$

(iv) We have to find out the value of $\bar{Y}$

where $\bar{Y} = \Sigma \hat{Y}_i$ ...

$$= (-0.25)^2 + (0.5)^2 + (0.75)^2 + 0.75 \ldots$$

In particular when $\bar{Y} = \bar{Y} = 30$, $X_1 = 4$, $X_{21} = 10$ then

$$\hat{Y}_i = \hat{\beta}_i = 23.75 - 0.25 \times 4 + 5.5 \times 10 = 30.5$$

When $Y = Y_4 = 24$, ...

then $\hat{Y} = \hat{\beta}_4 = ...$

$e_4 = Y_4 - \hat{Y}_4 = 24 ...$

When $Y = Y_5 = 40$, $\hat{Y}_5 = ...$

then $\hat{Y} = \hat{\beta}_5 = ...$

$= 23.75 ... 2 + 64 ... 40 ...$

$e_5 = Y_5 - \hat{Y}_5 = 40 - 40.75 = -0.75$

v) Now we have to calculate the variances of the OLS estimators of the regression parameters, $\text{var}(\hat{\beta}_1)$, $\text{var}(\hat{\beta}_2)$ and $\text{var}(\hat{\beta}_3)$.

We know that $\text{var}(\hat{\beta}) = \dfrac{\sigma_u^2 \Sigma x_{2i}^2}{\Sigma x_{1i}^2 \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$ ... Here $\sigma_u^2$ is

not known and hence it is replaced by its unbiased estimator $\hat{\sigma}_u^2 = \dfrac{\Sigma e_i^2}{n}$

Here $\dfrac{\Sigma e_i^2}{n-1} = \dfrac{1.5}{5-3} = \dfrac{1.5}{2} = 0.75$

$\text{var}(\hat{\beta}_1) = \dfrac{\left(\dfrac{\Sigma e_i^2}{n-1}\right)\Sigma x_{2i}^2}{\Sigma x_{1i}^2 \Sigma x_{2i}^2 - (\Sigma x_1 x_{2i})^2} = \dfrac{0.75 \times 10}{16 \times 0 - (2)^2}$

$= \dfrac{7.5}{16} = 0.4687 \qquad \text{var}(\hat{\beta}_1) = 0.4687$

Similarly, $\operatorname{var}(\beta_2) = \ldots$ ... ... ...

and $\beta_{\ldots} = \ldots$

and $\operatorname{var}(\beta_{\ldots} = \ldots$ ... $\operatorname{var} \ldots \ldots$ ... $\beta$ ... ... $\operatorname{var} \ldots$

We put $\sigma_u = \ldots \ldots$ ... ... ... $\hat{\beta}_u$ ...

$\operatorname{var}(\beta) = 0.6663^{\ldots}$ $\operatorname{var}(\beta) = 0.75$ and $\operatorname{cov}(\beta, \beta) = \ldots$ ... obtained form ...

Thus, we have $\operatorname{var}(\sigma_u \ldots \ldots \hat{\beta}) = 0.6663^{\ldots} \ldots \hat{\sigma} = \ldots$ ... $0.6663$

$0.75 \ldots \ldots 0.75$

$\operatorname{var}(\beta_0) = 3(0.75)$

(vi) We have to find out the value of $\operatorname{cov}(\beta_1 \beta_2)$

Here $\operatorname{cov}(\beta_1, \beta) = \ldots \ldots$ ... Here $\sigma_u$ is not known and hence

replaced by the unbiased estimator $\sigma_u = \dfrac{\ldots}{n-1} = 0.75$

$\operatorname{cov}(\beta_1, \beta) = \dfrac{0.75 \times 2}{n - 10 + 23} = \dfrac{9}{16} = 0.5625$

$\operatorname{cov}(\beta_1, \beta) = 0.5625$

## 3.3 Properties of OLS Estimator Vector β

Let $Y_i = \beta_0 + \beta_1 X_i + \beta_2 t + \ldots + \beta_K X_{Ki} + u_i$ be a General Linear Regression model

In vector-matrix form the model takes the form $Y = X\beta + u$

where $Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_{n \times 1} \end{bmatrix}$ ... $\begin{bmatrix} 1 & X & Y_{2i} & Y_{Ki} \\ 1 & X_i & Y_{\ldots} & Y_{Ki} \\ \vdots & & & \\ X_{1n} & X_{2n} & Y_{Kn} \end{bmatrix}_{n \times K}$

$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix}_{K \times 1}$ and $u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}_{n \times 1}$

Property 2 : $\hat{\beta}$ is an unbiased estimator of $\beta$

Proof ...

... $(X'X)^{-1}X'u$ ...

where ...

... as well as $X'X$

...

Again $E[\hat{\beta}] = \beta + E[(X'X)^{-1}X'u]$

$= \beta + (X'X)^{-1}X'E(u) = E(\hat{\beta})$   [ $E(u) = 0$ ]

$\hat{\beta} = \beta$ where ...

This shows that the OLS estimator of $\beta$ is unbiased estimator ...

$$
E(\hat{\beta}) = 
\begin{bmatrix}
E(\hat{\beta}_1) \\
E(\hat{\beta}_1) \\
E(\hat{\beta}_2) \\
\vdots \\
E(\hat{\beta}_K)
\end{bmatrix}
=
\begin{bmatrix}
\beta_1 \\
\beta_1 \\
\beta_2 \\
\vdots \\
\beta_K
\end{bmatrix}
$$

This implies that OLS estimator of each parameter is an unbiased estimator.
In terms of a linear regression model with two explanatory variables we have

$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$   $i = 1, 2, \ldots N$

In this model we have $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$   $\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}$

and $E(\hat{\beta}) = \begin{bmatrix} E(\hat{\beta}_0) \\ E(\hat{\beta}_1) \\ E(\hat{\beta}_2) \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}_{3 \times 1}$

$E(\hat{\beta}) = \beta$ where $E(\hat{\beta}_0) = \beta_0$, $E(\hat{\beta}_1) = \beta_1$, $E(\hat{\beta}_2) = \beta_2$

**Property 2** The Dispersion matrix or variance-covariance matrix of $\hat{\beta}$ is given by $\sigma_u^2(X'X)^{-1}$

**Proof** By definition dispersion matrix or variance-covariance matrix of $\hat{\beta}$ is given by $D(\hat{\beta}) = E[\hat{\beta} - E(\hat{\beta})][\hat{\beta} - E(\hat{\beta})]'$ where $E(\hat{\beta}) = \beta$

$$
\begin{bmatrix}
\beta_0 & \beta_1 \\
\beta_1 & \beta_1 ... \beta_1
\end{bmatrix}
$$

$$
\begin{bmatrix}
var(\hat{\beta}_0) & cov(\hat{\beta}_0,\hat{\beta}_1) & cov(\hat{\beta}_K,\hat{\beta}_0) \\
cov(\hat{\beta}_0,\hat{\beta}_1) & var\hat{\beta}_1 & cov(\hat{\beta}_K,\hat{\beta}_1) \\
cov(\hat{\beta}_K,\hat{\beta}_K) & cov(\hat{\beta}_1,\hat{\beta}_K) & var(\hat{\beta}_K)
\end{bmatrix}
$$

Here the diagonal terms are variances and non-diagonal terms are covariances. It is also called variance covariance matrix.

$$D(\hat{\beta}) = E[\hat{\beta} - \beta][\hat{\beta} - \beta]'$$

$$= E[(X'X)^{-1}X'u][(X'X)^{-1}X'u]'$$

$$= E[(X'X)^{-1}X'uu'X(X'X)^{-1}]$$

$$= (X'X)^{-1}X'\sigma_u^2 I_n X(X'X)^{-1}$$

$$= \sigma_u^2(X'X)^{-1}X'X(X'X)^{-1}$$

$$= \sigma_u^2(X'X)^{-1}$$

$$D(\hat{\beta}) = \sigma_u^2(X'X)^{-1}$$

Since $\hat{\beta} = \beta + (X'X)^{-1}X'u$, $\hat{\beta} - \beta = (X'X)^{-1}X'u$

Since $D(u) = E[u - E(u)][u - E(u)]' = E(uu')$ as $E(u) = 0 = \sigma_u^2 I_n$ (See 2.1 (u).

where $I_{n+1}$ = Identity matrix of order

Proceeding in the same way we can also derive the result $D(\hat{\beta}) = \sigma_u^2(X'X)^{-1}$ for a regression model with two explanatory variables

i.e. $Y = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$, $i = 1, 2, 3, ..., n$

**Property 3.** $j$th element of $\hat{\beta}$ is the best linear unbiased estimator of the $j$th element of $\beta$. Alternatively, $\hat{\beta}$ is the Best Linear Unbiased Estimator (BLUE) of $\beta$

**Proof** Since, $\hat{\beta} = (X'X)^{-1}X'Y$

Now we have to find out the conditions under which $\beta_0^*$ is an unbiased estimator of $\beta_0$.

Now, $\beta_0^* = C'Y = C'(X\beta + u)$ since $Y = X\beta + u$

$$= C'X\beta + C'u$$

$$E[\beta_0^*] = E[C'X\beta] + E[C'u]$$

$$= C'X\beta + 0 = C'X\beta \qquad [\because E(u) = 0_{n}.$$

$$E(\beta_0^*) = C'X\beta$$

Now $E(\beta_0^*) = \beta_0$ if $C'X = e_1'$

$$\text{as } e_1'\beta = [1 \quad 0 \qquad 0]\begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix} = \beta_0$$

This shows that $\beta_0^*$ is an unbiased estimator of $\beta_0$. The condition for $\beta_1^*$ to be an unbiased estimator of $\beta_0$ is given by $C'X = e_1$

Again, $\beta_0^* = C'X\beta + C'u$

$$= e_1'\beta + C'u = \beta_0 + [C_1 \quad C_2 \qquad C_n]\begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}$$

or, $\beta_0^* - \beta_0 \sum$ ... $\beta^* \beta_0 \sum$ ...

Now $\text{var}(\beta_0^*) = E[\beta_0^* - \beta_0]^2$ ... $E(\hat{\beta}) = \beta_0$

$$E \sum c ... = \sigma_u \sum c ...$$

$$\text{var} \beta_0^* = \sigma_u \sum c$$

Now we have to minimise $\text{var}(\beta_0^*)$ subject to the condition that through the choice of the vector $C$. In other words we have to minimise $\sigma_u^2 \sum c$ subject to the condition $C'X = e_1'$

For the sake of simplicity we put $\sigma_u^2 = 1$

The Lagrangian is given by,

$$L = \frac{1}{2} \sum_{i=1}^{n} C^2 - [C'X - e_1]\lambda$$

where $\lambda = \begin{bmatrix} \lambda_0 \\ \lambda \end{bmatrix}$ is the vector of Lagrangian multipliers

$$L = \frac{1}{2} C'C - [C'X - e_1]\lambda$$

Now differentiating it with respect to $C$ we get

$\frac{\partial L}{\partial \lambda} = C - X\lambda = 0_{n\times1}$ a null column vector where $C'C$ ... $C$

$\frac{\partial}{\partial C}(C'C) = 2C$ and $\frac{\partial L}{\partial \lambda} = C'X - e_1 = 0_{1\times k}$ ...

a null row vector ... $C = C_n' \sum_{i=1}^{n}$

$$\frac{\partial L}{\partial C} = C' - X\lambda = 0_{n\times1}$$

or $C = X\lambda$ i.e. $C = \lambda X$ and $C'X = \lambda X'X$

Again $\frac{\partial L}{\partial \lambda} = C'X - e_1 = 0_{1\times(k+1)}$ ... $e_1 = C'X$

$e_1 = C'X = \lambda X'X$ or, $\lambda X'X = e_1$

or, $\lambda' = e_1(X'X)^{-1}$ ... $C = \lambda X = e_1(X'X)^{-1}X$

so $\hat{\beta}$ ...

$$\beta_0$$
$$\beta$$
$$\beta_k$$

Now, it is proved that under the condition that $\beta^*$ is an unbiased estimator, the variance of $\beta^*$ i.e. var($\beta^*$) is minimum when $\beta^* = \beta$.

Applying the same mathematical technique it can be proved that

$$\text{var}(\beta_1^*) \text{ is minimum when } \beta_1^* = \beta_1$$

$$\text{var}(\beta_2^*) \text{ is minimum when } \beta_2^* = \beta_2$$

$$\text{var}(\beta_k^*) \text{ is minimum when } \beta_k^* = \beta_k$$

Hence it is proved that OLS estimators are the best linear unbiased estimators of the regression parameters i.e. $\hat{\beta}$ is the BLUE of $\beta$. This is known as the GAUSS-MARKOV THEOREM.

In case of three variable Linear regression model $y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + u_i$, $i = 1, 2, \ldots, n$) we can prove that $\hat{\beta}$ is the BLUE of $\beta$ where

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \text{ and } \hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{bmatrix}$$

**Property 4** Unbiased estimator of $\sigma_u^2$ is $\dfrac{e'e}{n - (K+1)} = \dfrac{\sum_1^n e_i^2}{n - (K+1)}$ where

$K$ = number of explanatory variables and $(K + 1)$ = number of parameters including the constant intercept term.

**Proof** Since $y = X\beta + u$, $Y = X\hat{\beta}$ and $y = \hat{y} + e$

or, $e = y - \hat{y}$ or $e = Y - X\hat{\beta}$    $e = Y - X\hat{\beta}$

Now $e'e = (Y - X\hat{\beta})'(Y - X\hat{\beta})$

$= [Y - X(X'X)^{-1}X'Y]'[Y - X(X'X)^{-1}X'Y]$    since $\hat{\beta} = (X'X)^{-1}X'Y$

$= [Y' - Y'X(X'X)^{-1}X'][Y - X(X'X)^{-1}X'Y]$

$= Y'[I - X(X'X)^{-1}X'][I - X(X'X)^{-1}X']Y$

$= Y'MMY$ where $M = I - X(X'X)^{-1}X'$

$\varepsilon'\varepsilon = \Sigma(M() )\;M)$                                    $I = $ Identity matrix

$= (Y\beta + u\;M)(\beta + u)$                          where $M$ is idempotent mat.

$= (\beta'X + u)M(X\beta + u)$                           for which $M^2 = M$

$= \beta'X'MX\beta + u'MX\beta + \beta'X'Mu + u'Mu$

Now if we put $M = I - X(X'X)^{-1}X'$, then

$\beta'X'MX\beta + u'MX\beta + \beta'X'Mu = 0$

and hence $\varepsilon'\varepsilon = u'Mu = u'[I - X(X'X)^{-1}X']u$

$$= u'[u - X(X'X)^{-1}X'u] = \sum_{i=1}^{n} u_i^2 - u'X(X'X)^{-1}X'u$$

$E\;\varepsilon'\varepsilon = \Sigma E(u_i^2) - E[u'u\;\text{trace}\;X(X'X)^{-1}X'] \quad E(u_i^2) = \sigma_u^2 \text{ and } u'u = \sum_{i}^{n} u_i^2$

$= n\sigma_u^2 - \sigma_u^2 (K-1) [\text{Here trace}\;X(X'X)^{-1}X' = K - 1 + 1]$

$E(\varepsilon'\varepsilon) = \sigma_u^2[n - (K+1)]$

or, $E\left[\dfrac{\varepsilon'\varepsilon}{n - K - 1}\right] = \sigma_u^2$ or, $E\left[\dfrac{\Sigma e_i^2}{n - (K+1)}\right] = \sigma_u^2$

This proves that $\dfrac{\varepsilon'\varepsilon}{n - (K+1)} = \dfrac{\Sigma e_i^2}{n - (K+1)}$ is the unbiased estimator of $\sigma_u^2$

In particular for a linear regression model with two explanatory variables

$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$

we have $E\left[\dfrac{\Sigma e_i^2}{n - 3}\right] = \sigma_u^2$ or, $E[\hat\sigma_u^2] = \sigma_u^2$ where $\hat\sigma_u^2 = \dfrac{\Sigma e_i^2}{n - 3}$ and $K = 2$

## 3.4. MLE of $\beta$ and $\sigma_u^2$ in the Multiple Regression Model

Let $Y = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_K X_{Ki} + u_i$ for $i = 2 \cdots n$ be the equation
the general linear regression model. In vector matrix form the set of $n$ equations can
be written as, $Y = X\beta + u$

where $Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n\times 1}$    $X = \begin{bmatrix} 1 & X_{11} & X_{21} & X_{K1} \\ 1 & X_{12} & X_{22} & X_{K2} \\ \vdots & & & \\ 1 & X_{1n} & X_{2n} & X_{Kn} \end{bmatrix}_{n\times K}$

$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix}_{(K+1)\times 1}$     $u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}_{n\times 1}$

... We also assume that ... independent

... the ... probability density functions of the ... given ... and ... the p.d.f. of $u_i$ for $i = 1, 2, \ldots, n$

Since $u_1, u_2, \ldots, u_n$ are independent ... $f(u_1, u_2, \ldots, u_n) = \prod f_i(u_i)$

Since $u_i \sim N(0, \sigma_u)$ then

$$f(u) = \frac{1}{\sqrt{2\pi}\,\sigma_u}\, e^{-\frac{u_i^2}{\sigma_u^2}} \quad (-\infty < u < \infty)$$

So, $f(u_1, u_2, \ldots, u_n) = \prod \frac{1}{\sqrt{2\pi}\,\sigma_u}\, e^{\cdot} \ldots e^{-\frac{\sum u_i^2}{\sigma_u^2}}$

This is the likelihood function of the parameters $\beta_0, \beta_1, \ldots, \beta_K$ and $\sigma_u$ and is denoted by

$$L(\beta', \sigma_u) = \frac{1}{(\sqrt{2\pi})^n \sigma_u^n}\, e^{-\frac{\sum u_i^2}{\sigma_u^2}} \quad \text{where } \beta' = [\beta_0, \beta_1, \ldots, \beta_K]_{(K+1)\times 1}$$

Now $\log L = -\frac{n}{2}\log 2\pi - n\log\sigma_u - \frac{1}{2\sigma_u^2}\sum_{i=1}^{n} u_i^2$

Now to obtain the MLE of the parameters Log $L$ is to be maximised with respect to the parameters

(a) To obtain the MLE of $\beta$ we have to maximise Log $L$ with respect to $\beta$ which is equivalent to minimization of $\sum_{i=1}^{n} u_i^2$ with respect to $\beta$.

So, we have to minimise $\sum u_i^2$ through the choice of $\beta$

Since $u = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_N \end{bmatrix}_{N\times 1} \quad u' = [u_1\ u_2\ \ldots\ u_n]_{1\times n}$

$$u'u = \sum_{i=1}^{n} u_i^2$$

160-9

It should be noted that MLE of $\sigma_u^2$ is not an unbiased estimator of , but a consistent or asymptotically unbiased estimator of $\sigma_u^2$

Since MLE of $\sigma_u$  $\sum_{i} e_i^2$  $n$

Now  $\sum e_i^2 / n = \dfrac{n-(K+1)}{n} \cdot \dfrac{\sum_{i=1}^{n} e_i^2}{n-(K-1)} = \left(1 - \dfrac{K+1}{n}\right) \cdot \dfrac{\sum_{i=1}^{n} e_i^2}{n-K}$

or  $\sum_{0}^{n} e_i^2 / n = \left(\dfrac{K+1}{n}\right) \dfrac{\sum_{i} e_i^2}{n-(K+1)}$

As $N \to \infty$, $\dfrac{(K+1)}{n} \to 0$ and hence

$\sum_{i=1}^{N} \sigma_i^2 / n \to \sum_{i}^{n} e_i^2 / n-(K+1)$ where $\dfrac{\sum_{i=1}^{n} e_i^2}{n-(K+1)}$ is an unbiased estimator of $\sigma_u^2$

This proves that MLE of $\sigma_u^2$ i.e $\sum_{i}^{n} e_i^2 / n$ is an asymptotically unbiased or consistent estimator of $\sigma_u^2$. In a three variable (with two explanatory variables, i.e., when $K = 2$) linear regression model we have MLE of $\sigma_u^2 = \sum_{i}^{n} e_i^2 / n$ and unbiased estimator of $\sigma_u^2 = \sum_{i=1}^{n} e_i^2 / (n-3)$

## 3.5 Expression of Multiple Correlation Coefficient in the General Linear Regression Model

Let $\quad \mu_0 \quad \beta_1 X_{1i} \quad \beta_2 X_{2i} \quad \beta_3 X_{3i} \quad + \beta_k X_k + u_i$ be the regression

where $\quad 2 \quad n$

Here $Y$ is regressed on $X_1 \ X_2 \ X_3 \quad X_k$. So the multiple correlation coeff

denoted by the symbol,

$$R_{\cdot} \quad X_k = \frac{\operatorname{var} Y}{\Sigma} \quad \frac{\sum Y \ \bar{Y}_i \ n}{\sum Y \ \bar{Y} \ u} \qquad \text{Since } Y = Y$$

$$= \frac{\sum \hat{Y}_i^2 - n\bar{Y}^2}{\sum Y_i^2 - n\bar{Y}^2} = \frac{\hat{Y}'\hat{Y} \ n^2}{Y'Y - n\bar{Y}^2}$$

$$= \frac{\beta \ Y|| \ n\bar{Y}^2}{Y'Y \ n\bar{Y}^2} \qquad \text{where } \hat{Y} = X\hat{\beta} \text{ and } \hat{Y} = Y$$

$$= \frac{\hat{\beta} \ X'Y\hat{\beta} \ n^2}{Y \ Y - n\bar{Y}^2} \quad \frac{\beta \ (X'Y - n\bar{Y}^2)}{Y'Y - n^2} \qquad \begin{array}{l} Y \ x_1 Y \ Y \quad Y_{n \ k \ n} \\ \hat{\beta} = (X' \quad n \ Y \\ \hat{\beta} \ )X \ X'Y \end{array}$$

where $\hat{\beta} = X \ X'Y \quad (X'X)^{-1} X'Y$

## 3.6 The Multiple Coefficient of Determination $R^2$ and the Multiple Coefficient of Correlation in the Three-Variable Linear Regression Model

In the two-variable case we have seen that $R^2$ (or $r^2$) measures the goodness of fit of the regression equation $(Y = \alpha \quad \beta X + u, \quad i \ 2 \quad n)$ that is it gives the proportion or percentage of the total variation in the dependent variable $Y$ explained by the single explanatory variable $X$. This notion of $R^2$ can be easily extended to regression models containing more than two variables. Thus in the three-variable model we would like to know the proportion of the variation in $Y$ explained by the variables $X$ and $X_2$ jointly.

This quantity that gives this information is known as the multiple coefficient of determination and is denoted by $R^2_{Y \ X_2 \ X_3}$ or simply $R^2$ and conceptually it is similar to $r^2$.

The estimated three-variable regression line $(Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \ i = 1, 2$

$n)$ is given by $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} \ \hat{\beta}_2 X_{2i}$ where $\hat{Y}_i$ is the estimated value of $Y$. Both the fitted regression line and is an estimator of true $E(Y, X_{2i} \ Y_{3i})$ and $\bar{Y} = \beta_0 \ \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$.

Taking deviations from means we have

$$\hat{Y}_i \ \bar{Y} \quad \beta_0 + \beta_1 X_{1i} + \hat{\beta}_2 Y \quad \beta_0 \ \hat{\beta}_1 \bar{X}_1 \ \hat{\beta}_2 \bar{X}_2$$

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \ldots$$

$$\hat{y}_i = \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}$$

where $y = Y - \bar{Y}$, $x = X - \bar{X}$ and $e = r$

Now errors of estimate $e_i = y_i - \hat{y}_i = y_i - (\hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i})$

$y_i = \hat{y}_i + e_i$

Now $\Sigma y^2 = \Sigma \hat{y}^2 + \Sigma e^2 + 2\Sigma \hat{y} e$

or $\Sigma y^2 = \Sigma \hat{y}^2 + \Sigma e^2$   $\Sigma \hat{y} e = 0$

i.e. TSS = ESS + RSS

when   TSS = Total Sum of Squares

ESS = Explained Sum of Squares

RSS = Residual Sum of Squares

Now by definition, $R^2_{Y X_1 X_2} = R^2 = \dfrac{\Sigma \hat{y}_i^2}{\Sigma y^2} = \dfrac{ESS}{TSS}$

$$= \frac{\Sigma(\hat{Y}_i - \bar{Y})^2}{\Sigma(Y - \bar{Y})^2} = \frac{\Sigma y^2 - \Sigma e^2}{\Sigma y_i^2}$$

Since

$\Sigma y_i^2 = \Sigma \hat{y}^2 + \Sigma e^2$

$$= 1 - \frac{\Sigma e^2}{\Sigma y^2} = 1 - \frac{RSS}{TSS}$$

$\Sigma \hat{y}_i^2 = \Sigma y_i^2 - \Sigma e^2$

Since $e_i = y - \hat{y}_i$,   $\Sigma e_i^2 = \Sigma e_i (y_i - \hat{y}_i)$

$= \Sigma e_i(y_i - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i})$ as $\hat{y}_i = \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}$

$= \Sigma e_i y_i - \hat{\beta}_1 \Sigma x_{1i} e_i - \hat{\beta}_2 \Sigma x_{2i} e_i$

$= \Sigma e_i y_i$

$= \Sigma(y_i - \hat{y}_i) y_i$

$= \Sigma y_i (y_i - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i})$

$= \Sigma y_i^2 - \hat{\beta}_1 \Sigma x_{1i} y_i - \hat{\beta}_2 \Sigma x_{2i} y_i$

$\Sigma e_i^2 = \Sigma y_i^2 - \hat{\beta}_1 \Sigma x_{1i} y_i - \hat{\beta}_2 \Sigma x_{2i} y_i$

Since $\Sigma x_{1i} e_i = 0$, $\Sigma x_{2i} e_i = 0$

where $\Sigma x_{1i}(y_i - \hat{\beta}_1 x_{1i} - \hat{\beta}_2 x_{2i})$

$= \Sigma x_{1i} y_i - \hat{\beta}_1 \Sigma x_{1i}^2 - \hat{\beta}_2 \Sigma x_{1i} x_{2i}$

$= 0$

which follows from the first normal equation.

Similarly $\Sigma x_{2i} e_i = 0$

i.e., $RSS = \Sigma y_i^2 - \hat{\beta}_1 \Sigma x_{1i} y_i - \hat{\beta}_2 \Sigma x_{2i} y_i$

Since $\Sigma y^2 = TSS$

$ESS = \hat{\beta}_1 \Sigma x_{1i} y_i + \hat{\beta}_2 \Sigma x_{2i} y_i$

$$R^2_{Y X_1 X_2} = \frac{ESS}{TSS} = \frac{\hat{\beta}_1 \Sigma x_{1i} y_i + \hat{\beta}_2 \Sigma x_{2i} y_i}{\Sigma y^2}$$

**Example 3.3** ...

(i) find the value of $R$

(ii) Find the fitted regression equation

(b) Following Example 3.2,

(i) find the value of $R$

... the ... give an equation and interpret

**Solution**

(a) From the above table of Example ... we get $\hat{\beta}$ ...

(b) ... we obtained ... and $\hat{\beta}$ ...

We know that $R$ ... $= \dfrac{\hat{\beta}_1 \Sigma x_{1i} y_i + \hat{\beta}_2 \Sigma x_{2i} y_i}{\ldots}$

... 

$$R^2 = 0.7 ...$$

... estimated fitted regression equation is

$$\hat{y} = \ldots$$

... and $R = .7$ ...

(b) ... the values obtained of Example 3.2

we get ... $n_1^2 \ldots = 42$, $y_i = \ldots$

Further we obtained $\hat{\beta}_1 \ldots$ ... $\hat{\beta}_2 = 0.25$ and $\hat{\beta}_3 = 5.5$

we thus get $R^2 = \hat{\beta} \ldots$ $+ \ldots$

$$= \frac{4.25 \, n_i \ldots}{27} \ldots + \frac{5}{2} \, 9i \ldots = \frac{95}{\ldots} = 0.994$$

$$R = 0.994$$

(ii) Thus the estimated regression equation is

$$\hat{y}_i = 23.75 + 0.25x \ldots + 5.5x_i, \quad R = 0.994$$

This equation suggests that years of experience with the firm is far more important than years of education (which actually has a certain negative sign). This equation says that we can predict that one more year of experience, after allowing for years of education (for holding it constant), results in an annual increase in salary of $5500. This means that if we consider the persons with the same level of education, the one with one more year of experience can be expected to have a higher salary of $5500. Similarly, if we consider two persons with the same experience, the one with an education of one more year can be expected to have a lower annual salary of $ .50.

Here $R^2 = 0.995$ implies that out of 100% variation in salary of the employee, 99% variation can be explained by the two explanatory variables $X_2$ and $X_3$ jointly.

## 7.7 $R^2$ and the Adjusted $R^2$

[faded text] ... will not decrease $R^2$

[illegible lines]

Now TSS ... is the product of ...

... ... ...

The RSS $\Sigma_i$ ... however depends on ...

[illegible] $K$ ... part but as the number of ... is likely to decrease
(at least will not increase) hence $R^2$ will ... [illegible] comparing ...
regression models with the same dependent variable ...
variables, one should be very wary of choosing the model ... $R^2$

To compare two $R^2$ terms one must take into account the number of variables
present in the model. This can be done readily if we consider a different ... [illegible]
of determination,

$$\bar{R}^2 = 1 - \frac{\Sigma e_i^2}{\Sigma y_i^2} \frac{(n - (K+1))}{(n-1)} \qquad \text{where } K = \text{number of explanatory variables and } K$$

= number of parameters in the model including the intercept term. In a three variable
(with two explanatory variables) linear regression model, $K = 2$ ... $K$ ... $n$ ...

The $R^2$ thus defined is known as adjusted $R^2$ denoted by $\bar{R}^2$. The term adjusted
means adjusted for the degrees of freedom (d.f.) associated with ... sum of squares
of $\Sigma e_i^2$ and $\Sigma y_i^2$

RSS $= \Sigma e_i^2$ has $n - (K+1)$ degrees of freedom in a model involving (K ... )
parameters, including the intercept term and TSS $= \Sigma y_i^2$ has $n - 1$ degrees of freedom

Thus the adjusted $R^2$ can also be written as

$$\bar{R}^2 = \frac{\hat{\sigma}_u^2}{S_y^2}$$

It is thus clear that
$\bar{R}^2$ and $R^2$ are
related and we can
express the relation
as follows

where $\hat{\sigma}_u^2 = \Sigma e_i^2 / (n - (K+1))$ is the residual variance and
unbiased estimator of true $\sigma_u^2$ and

$$S_y^2 = \frac{1}{n-1} \Sigma (y_i - \bar{y})^2 = \frac{1}{n-1} \Sigma y_i^2 \qquad \text{sample variance of } y$$

$$\Sigma y_i^2 = (n-1)S_y^2 \text{ and } \Sigma y_i^2 / (n-1) = S_y^2$$

...

### Note Comparing Two $R^2$ values

It is important to note that in comparing two models on the basis of the coefficient of determination, whether adjusted or not, the sample size $n$ and the dependent variable must be the same, the explanatory variables may take any form. Thus for the models

$$\log Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i \qquad \text{(A)}$$

and

$$Y_i = \alpha_0 + \alpha_1 X_{1i} + \alpha_2 X_{2i} + u_i \qquad \text{(B)}$$

the computed $R^2$ terms cannot be compared. The reason is that by definition $R^2$ measures the proportion of the variation in the dependent variable as summarized by the explanatory variables. Therefore in equation (A) $R^2$ measures the proportion of the variation in $\log Y$ explained by $X_1$ and $X_2$ whereas in equation (B) it measures the proportion of the variation in $Y$ and hence the two are not the same thing. A change in $\log Y$ gives a relative or proportional change in $Y$ whereas a change in $Y$ gives an absolute change. Therefore $\text{var}(\log Y)$ and $\text{var}(Y)$ is not equal to ... Thus the two coefficients of determination are not the same.

**Example 3.4.** a) Following Example 3.1 find the value of Adjusted $R^2$
b) Following Example 3.2 find the value of Adjusted $R^2$.

**Solution** (a) We know that in a three-variable linear regression model, adjusted $R^2$ is denoted by

$$\bar{R}^2 = 1 - (1 - R^2)\frac{n-1}{n-3}$$ Here we see that following data of Example ... $R^2 = $ ...

$n = 5$

... the value of adjusted $R^2 = \bar{R}^2 = 0.13$

... 

... allowing Example ... the ... $\frac{R}{n}$ ... $\frac{1}{3}$

$n = 3$ ... and $n = 3$

$N$ ... $R$ ... $\frac{n}{n}$ ...

$1.013 = 0.988$

Adjusted $R$ ... $0.988$ which is smaller than unadjusted $R$ ...

The value of adjusted $R$ ... $\bar{R}$ ... can also be obtained ...

$$\bar{R}^2 = 1 - \frac{\Sigma e^2 / (n - (K+1))}{\Sigma_{y^2} / (n-1)}$$

From Example 3.2 we have obtained

$$\Sigma e^2 = 15, \quad \Sigma_{y^2} = 272, \quad n = 5, \quad K = 2, \quad K + 1 = 3$$

$$\bar{R}^2 = 1 - \frac{\frac{15}{5-2}}{\frac{272}{5-1}} = \frac{0.75}{68} = 1 - 0.0110 = 0.988$$

$$\bar{R}^2 = 0.988$$

## 3.8. Partial Correlation Coefficients and the Coefficient of Partial Determination

In the simple correlation analysis the coefficient of correlation $r$ is used as a measure of the degree of linear association between two variables $X$ and $Y$

$$[Y = \alpha + \beta X_i + u_i, \quad i = 1, 2, \dots, n]$$

For three variable Linear regression model $[Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad i = 1, 2, \dots, n]$ we can compute three correlation coefficients $r_{YX_1} = r_{12}$ (correlation coefficient between $Y$ and $X_1$), $r_{YX_2} = r_{13}$ (correlation coefficient between $Y$ and $X_2$) and $r_{X_1 X_2} = r_{23}$ (correlation coefficient between $X_1$ and $X_2$)

These correlation coefficients are called gross or simple correlation coefficients or correlation coefficients of zero order and computed by the formula

$$r_{XY} = \frac{\operatorname{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\frac{1}{n}\Sigma(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\frac{1}{n}\Sigma(X_i - \bar{X})^2}\sqrt{\frac{1}{n}\Sigma(Y_i - \bar{Y})^2}} = \frac{\Sigma x_i y_i}{\sqrt{\Sigma x^2}\sqrt{\Sigma y^2}}$$

where $x_i = X_i - \bar{X}$ and $y_i = Y_i - \bar{Y}$

## Partial correlation coefficient

Partial correlation between $x$ and $y$

Partial correlation between $y$ and $z$

The partial correlation defined above are called

From the formula

$$r = \frac{}{\sqrt{\quad}} \qquad \text{for example}$$

... and ...
$X_1$ are uncorrelated

The ... coefficient of partial determination is ... be interpreted as the proportion of the variation in $Y$ not explained by the variable $X_2$ but has been explained by the inclusion of $X_1$ into the model. Conceptually it is another ... coefficient of determination ...

... using ... are found ... and partial correlation coefficients

$$R = \cdots$$

$$R^2 = \cdots$$

$$R = \cdots$$

... has been pointed out earlier that $R^2$ will not decrease ... explanatory variable is introduced into the model which ... can be seen clearly that the equation ...
$R^2 = \eta_2^2 + \cdots$ This equation states that the proportion ... the variation in $Y$ explained by $X_1$ and $X_2$ jointly is the sum of two parts: the part explained by $X_2$ and the ... $\eta_2^2$ and the part not explained by ... times the proportion that is explained by $X_2$ after holding the influence of $X_1$ constant. Now $R^2 \geq \eta_2^2$ so long as ... As worst, $\eta_{12}^2$ will be zero, in which case $R^2 = \eta_2^2$.

**Example 3.5.** Following Example 3.1
(i) Find the values of $r_{12}$, $r_{13}$ and $r_{23}$.
(ii) Find the values of the partial regression coefficient $r_{123}$.
(iii) Find the value of $R^2$ in terms of $r_{12}$, $r_{13}$ and $r_{23}$.

**Solution** From the calculation table of Example 3.1 we have the following values
$\Sigma x_1^2 = 40$, $\Sigma x_2^2 = 23.20$, $\Sigma y_i^2 = 24$, $\Sigma x_1 x_2 = 17$, $\Sigma x_1 y_i = 20$ and $\Sigma x_2 y_i = 3$
where $x_{1i} = X_{1i} - \bar{X}_1$, $x_{2i} = X_{2i} - \bar{X}_2$, $Y_i$ and $i = 1, \ldots$

(i) Now by using the formula of $r_{12}$, $r_{13}$ and $r_{23}$ and putting the required values we can get the value of $r_{12}$, $r_{13}$ and $r_{23}$.

By definition, $r_{12} = r_{Y} = \dfrac{\text{cov } x_1 y_1}{\sigma_Y \; \sigma_Y} = \dfrac{\frac{1}{n}\Sigma(x_1 - \bar{x}_1)(y_1 - \bar{y}_1)}{\sqrt{\frac{1}{n}\Sigma x_1 x_1 - \bar{x}^2} \sqrt{\frac{1}{n}\Sigma() - }}$

$$= \dfrac{\frac{1}{n}\Sigma x_1 y_1}{\sqrt{\frac{1}{n}\Sigma x_1^2 \cdot \frac{1}{n}\Sigma y_1^2}} = \dfrac{\Sigma x_1 y_1}{\sqrt{\Sigma x_1^2 \sqrt{\Sigma y_1^2}}} = \dfrac{20}{\sqrt{40 \times 24}} = 0.64$$

$$r_{12} = 0.64$$

Similarly, ...

and ...

$$r_{13} = 0.64, r_{13} = -0.13 \text{ and } r_{23} = 0.56$$

We know that ...

$$r_{12} = 0.50$$

(iii) We know that,

$$R^2 = \frac{r_{12}^2 + r_{13}^2 - 2 r_{12} r_{13} r_{23}}{1 - r_{23}^2}$$

$$= \frac{(0.64)^2 + (-0.13)^2 - 2 \times 0.64 \times (-0.13) \times 0.56}{1 - (0.56)^2}$$

$$= 0.76$$

$$R = 0.76$$

**Example 3.6.** Are the following data consistent ? Give reasons

(a) $r_{23} = 0.9$, $r_{31} = -0.2$, $r_{12} = 0.8$

(b) $r = 0.6$, $r = -0.9$, $r = 0.5$

**Solution** From the above data set we will find calculate the value

$$R^2 = \frac{r_{12}^2 + r_{13}^2 - 2 r_{12} r_{13} r_{23}}{1 - r_{23}^2}$$ and will verify whether $0 < R^2$ or not

(a) Here we see that

$$R^2 = \frac{r_{12}^2 + r_{13}^2 - 2 r_{12} r_{13} r_{23}}{1 - r_{23}^2}$$

$$= \frac{(0.8)^2 + (-0.2)^2 - 2 \times 0.8 \times (-0.2) \times 0.9}{1 - (0.9)^2} = 5$$

$R^2 = 5$ which is not possible as $0 < R^2$

Hence the given information are not consistent

(b) $R^2 = \frac{r_{12}^2 + r_{13}^2 - 2 r_{12} r_{13} r_{23}}{r_{23}^2}$ as $r_{31} = r_{13} = 0.5$

$$= \frac{(0.6)^2 + (-0.5)^2 - 2 \times 0.6 \times 0.5 \times 0.9}{1 + (0.9)^2}$$

**Example ...**

... $= 6042.200$, $\Sigma(X_{1i} - \bar{X}_1)^2 = 84855.096$

... $= 280\,000$, $\Sigma(Y_i - \bar{Y})(X_{1i} - \bar{X}_1) = 74778.346$

... and ...

**Solution** The given results are treated in a three variable linear regression model of the form $Y = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$, $i = 1, 2, \dots n$

where $\beta_i$ and ... are the partial regression coefficients. Assuming that ...

... and $Y$, ... we can obtain the Least squares estimates ... as follows

$$\hat{\beta}_1 = \frac{\Sigma x_{2i}^2 \, \Sigma x_{1i} y_i - \Sigma x_{1i} x_{2i} \, \Sigma x_{2i} y_i}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$$

$$= \frac{280\,000 \times 74778.346 - 4796.00 \times 4250.900}{84855.096 \times 280.00 - (4796.00)^2}$$

$$= \frac{20937936.88 - 20187316.40}{23759426.88 - 23001616} = \frac{550620.48}{757810.88} = 0.7265$$

$$\beta = 0.7265$$

Similarly, $\hat{\beta}_2 = \dfrac{\Sigma x_{1i}^2 \, \Sigma x_{2i} y_i - \Sigma x_{1i} x_{2i} \, \Sigma x_{1i} y_i}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$

$$= \frac{84855.096 \times 4250.900 - 4796.00 \times 74778.346}{84855.096 \times 280.000 - (4796.00)^2}$$

$$= \frac{2073580.17}{75780.88} = 2.7362$$

$$\hat{\beta}_2 = 2.7362$$

Now we have to find out $SE(\hat{\beta}_1) = \sqrt{\text{var}(\hat{\beta}_1)}$ and $SE(\hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_2)}$

We know that $\text{var}(\hat{\beta}_1) = \dfrac{\sigma_u^2 \Sigma x_{2i}^2}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$

and $\text{var}(\hat{\beta}_2) = \dfrac{\sigma_u^2 \Sigma x_{1i}^2}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$

$SE(\beta_1) = \sqrt{var\,\beta_1} \cdots 0.0540$

Similarly, var$\beta_1 = $

$SE(\hat{\beta}_1) = \sqrt{var(\beta_1)} = 0.8878$

Now we have to find out the value of $R^2$ and adjusted $R^2 = \bar{R}^2$

We know that $R^2 = \dfrac{ESS}{TSS} = 1 - \dfrac{RSS}{TSS} = 1 - \dfrac{\sum e_i^2}{\sum y_i^2}$

[since $\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2 \Rightarrow TSS = ESS + RSS$]

so $R$ and $\bar{R}^2$ are almost the same

## Confidence Intervals and Hypothesis Testing in a three variable Multiple Linear Regression Model

We construct a confidence interval ...

given by
$$\beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i, \quad i = 1, 2, \ldots n$$

Suppose $\hat{\beta}_0$, $\hat{\beta}_1$, $\hat{\beta}_2$ and $\sigma_u$ are the OLS estimators ...

respectively

We also know that

$$\hat{\beta}_1 \sim \beta_1, \text{var}(\hat{\beta}_1) \text{ where } E(\hat{\beta}_1) = \beta_1 \text{ and var}(\beta_1) = \underline{\qquad} \qquad \text{This is read}$$

as $\hat{\beta}_1$ is normally distributed with mean $\beta_1$ and variance var$(\hat{\beta}_1)$

Similarly $\hat{\beta}_2 \sim \beta_2, \text{var}(\hat{\beta}_2)$ where $E(\hat{\beta}_2) = \beta_2$

$$\text{and var}\beta_2 = \frac{\sigma_u^2}{n\sigma_{X_1}^2 \ \ }$$

$$\hat{\beta}_0 \sim \left[ \beta_0, \text{var}(\hat{\beta}_0) \right] \text{ where } E(\hat{\beta}_0) = \beta_0$$

$$\text{and var}(\hat{\beta}_0) = \frac{\sigma_u^2}{n} + \bar{X}_1^2 \text{ var}(\hat{\beta}_1) + 2\bar{X}_1\bar{X}_2 \text{ cov}(\hat{\beta}_1, \hat{\beta}_2) + \bar{X}_2^2 \text{ var}(\hat{\beta}_2)$$

and $u \sim N(0, \sigma_u^2)$ where $E(u) = 0$ and var $(u) = \sigma_u^2$

$$\text{and cov}(\hat{\beta}_1, \hat{\beta}_2) = \frac{\sigma_u^2 r_{X_1 X_2}}{n\sigma_{X_1}\sigma_{X_2}(1 - r_{X_1 X_2}^2)}$$

We are now interested in testing the following hypothesis

**Case 1** We want to test the null hypothesis $H_0 : \beta_0 = 0$ against the alternative hypothesis, either $H_1 : \beta_0 \neq 0$ or $H_1 : \beta_0 > 0$ or $H_1 : \beta_0 < 0$

Since $\hat{\beta}_0 \sim N[\beta_0, \text{var}(\hat{\beta}_0)]$

Now $t$ or $Z = \dfrac{\hat{\beta}_0 - \beta_0}{SE(\hat{\beta}_0)} \sim N(0,1)$

would be the appropriate test statistic where $SE(\hat{\beta}_0) = \sqrt{\text{var}\hat{\beta}_0}$

When $\sigma_u$ is unknown and $\hat{\sigma}$ replaces by its estimate ...

Since $E\left[\dfrac{\Sigma e^2}{n-3}\right] = \sigma_u$, then the appropriate test statistic would be ...

will follow a ... distribution with d.f. $= n - ...$ under $H_0$: $\beta_0 = ...$. The test statistic will be

$$t = \frac{\hat{\beta}_0 - \beta_0}{\sqrt{\text{var}\,\hat{\beta}_0}} = \frac{\hat{\beta}_0 - \beta_0}{SE\,\hat{\beta}_0}$$

## Nature of the Test

(i) For the alternative hypothesis $H_1$: $\beta_0 \neq 0$, the null hypothesis $H_0$ ... be accepted at $100\alpha\%$ level of significance if for the given sample ... $-t_{\alpha/2} < t < t_{\alpha/2}$ and will be rejected otherwise, i.e. when ... $\alpha$.

(ii) For the alternative hypothesis $H_1$: $\beta_0 > 0$, the null hypothesis ... $\beta_0 = 0$ ... will accepted if for the given sample $t \leq t_{\alpha}$ and will be rejected otherwise when $t > t_{\alpha}$.

(iii) For the alternative hypothesis $H_1$: $\beta_0 < 0$, the null hypothesis $H_0$ ... be accepted if for the given sample $t_{\alpha} > t$ and will be rejected otherwise i.e. when $t < -t_{\alpha}$. In each case $\alpha$ denotes the chosen level of significance. Usually $\alpha = 0.05$ or $0.01$.

## Confidence Interval for $\beta_0$:

As regards the problem of interval estimation of $\beta_0$ at $100\alpha\%$ level of significance, the confidence limits to $\beta_0$ would be given by

$$\hat{\beta}_0 \pm t_{\alpha/2} \cdot SE\,\hat{\beta}_0$$

i.e. $P\left[-t_{\alpha/2} < t < t_{\alpha/2}\right] = 1 - \alpha$

or $P\left[-t_{\alpha/2} < \dfrac{\hat{\beta}_0 - \beta_0}{SE(\hat{\beta}_0)} < t_{\alpha/2}\right] = 1 - \alpha$

or $P\left[\hat{\beta}_0 - t_{\alpha/2} \cdot SE(\hat{\beta}_0) \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2} \cdot SE(\hat{\beta}_0)\right] = 1 - \alpha$

Here $(1 - \alpha)$ is called the confidence coefficient.

**Case 2** We want to test the null hypothesis $H_0$: $\beta_1 = 0$ against the alternative hypothesis $H_1$: $\beta_1 \neq 0$ or $H_1$: $\beta_1 > 0$ or $H_1$: $\beta_1 < 0$

Since $\hat{\beta}_1 \sim N(\beta_1, \text{var}\,\hat{\beta}_1)$ where $E(\hat{\beta}_1) = \beta_1$

and $\text{var}(\hat{\beta}_1) = \dfrac{\sigma_u^2}{n\sigma_{x_1}^2(1 - r_{x_1 x_2}^2)}$

Now $t$ or $Z = \dfrac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \sim N(0,1)$ would be the appropriate test statistic where

... $H_0$: $\beta_1 = 0$, the test statistic would be

### Nature of the Test

(i) ... the null hypothesis $H_0$: $\beta = 0$ will be accepted ... for the given sample ...

(ii) ... the alternative hypothesis ($H_1$: $\beta$ ...) ... be accepted ... for the given sample ... when ...

(iii) for the alternative hypothesis $H_1$: $\beta$ ... accepted ... for the given sample ... and ... when ... In each case $\alpha$ (= 0.01 or 0.05) denotes the chosen level of significance.

### Confidence interval for $\beta_1$

As regards the problem of interval estimation ... the confidence limits to $\beta_1$ would be given by,

$$\hat{\beta}_1 \pm t_{\alpha/2,\,n-3} SE(\hat{\beta}_1)$$

i.e. $P[-t_{\alpha/2,\,n-3} \le t \le t_{\alpha/2,\,n-3}] = 1 - \alpha$

or, $P\left[-t_{\alpha/2,\,n} \le \dfrac{\hat{\beta}_1 - \beta_1}{SE(\hat{\beta}_1)} \le t_{\alpha/2,\,n-3}\right] = 1 - \alpha$

or, $P[\hat{\beta}_1 - t_{\alpha/2,\,n-3} SE(\hat{\beta}_1) \le \beta_1 \le \hat{\beta}_1 + t_{\alpha/2,\,n-3} SE(\hat{\beta}_1)]$

Here $(1-\alpha)$ is the confidence coefficient.

**Case 3** We want to test the null hypothesis $H_0$: $\beta_2$ ... against the alternative hypothesis, $H_1$: $\beta_2 \neq 0$ or $H_1$: $\beta_2 > 0$ or, $H_1$: $\beta_2 < 0$

Since $\hat{\beta}_2 \sim N[\beta_2, \text{var}(\hat{\beta}_2)]$ where $E(\hat{\beta}_2) = \beta_2$

and $\text{var}(\hat{\beta}_2) = \dfrac{\sigma_u^2}{n\sigma_{x_1}^2(1 - r_{1,x}^2)}$

Now $t$ or $Z = \dfrac{\hat{\beta}_2 - \beta_2}{SE(\hat{\beta}_2)} \sim N(0,1)$ would be the appropriate test statistic where

$SE(\hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_2)}$. When $\sigma_u^2$ is not known then it is replaced by its unbiased

estimator $\hat{\sigma}_u^2 = \Sigma e_i^2/(n-3)$ and the test statistic becomes $t = \dfrac{\hat{\beta}_2 - \beta_2}{SE(\hat{\beta}_2)} \sim t_{n-3}$

... 60 ... 0

Under $H_0 : \beta_2 = 0$ the test statistic would be

$$t = \frac{\hat{\beta}_2 - 0}{SE(\hat{\beta}_2)} = \frac{\hat{\beta}_2}{SE(\hat{\beta}_2)}$$

## Nature of the Test

(a) For the alternative hypothesis $H_1 : \beta_2 \neq 0$, the null hypothesis $H_0$ will be accepted for the given sample if $|t| \leq t_{\alpha/2, n-2}$ and will be rejected otherwise.

(b) For the alternative hypothesis $H_1 : \beta_2 > 0$, the null hypothesis $H_0$ will be accepted for the given sample if $t \leq t_{\alpha, n-2}$ and will be rejected otherwise when $t > t_{\alpha, n-2}$.

(b) For the alternative hypothesis $H_1 : \beta_2 < 0$, the null hypothesis $H_0$ will be accepted for the given sample if $t \geq -t_{\alpha, n-2}$ and will be rejected otherwise. [i.e. when $t < -t_{\alpha, n-2}$]

## Confidence Interval for $\beta_2$.

As regards the problem of interval estimation of $\beta_2$ at $100\alpha\%$ level of significance the confidence limits to $\beta_2$ would be given by

$$\hat{\beta}_2 \pm t_{\alpha/2, n-2} SE(\hat{\beta}_2)$$

or $P\left[ t \leq t_{\alpha/2, n-2} \right] = 1 - \alpha$

or $P\left[ \dfrac{\hat{\beta}_2 - \beta_2}{SE(\hat{\beta}_2)} \leq t_{\alpha/2, n-2} \right] = 1 - \alpha$

or $P\left[ \hat{\beta}_2 - t_{\alpha/2, n-2} SE(\hat{\beta}_2) \leq \beta_2 \leq \hat{\beta}_2 + t_{\alpha/2, n-2} SE(\hat{\beta}_2) \right] = 1 - \alpha$

where $(1 - \alpha)$ is the confidence coefficient.

**Case 4** We want to test the null hypothesis $H_0 : \beta_1 = \beta_2$ against the alternative hypothesis $H_1 : \beta_1 \neq \beta_2$ or $H_1 : \beta_1 > \beta_2$ or $H_1 : \beta_1 < \beta_2$

Since $\beta_1 - \beta_2 = \lambda, \hat{\beta}_1 - \hat{\beta}_2 = \hat{\lambda}$ and $\lambda = 0$

where $E(\hat{\beta}_1) = \beta_1, E(\hat{\beta}_2) = \beta_2, E(\hat{\beta}_1 - \hat{\beta}_2) = \beta_1 - \beta_2$

and $\text{var}(\hat{\beta}_1 - \hat{\beta}_2) = \text{var}(\hat{\beta}_1) + \text{var}(\hat{\beta}_2) - 2\text{cov}(\hat{\beta}_1, \hat{\beta}_2)$

$$= \sigma_u^2 \frac{1}{X^2} + \sigma_u^2 \frac{1}{\sum x^2} + \ldots$$

so $\text{Cov}(\hat{\beta}_1, \hat{\beta}_2) = \dfrac{\sigma_u^2 \bar{X}_2}{n \sum x^2} \ldots$

The appropriate test statistic would be given by

$$t \text{ or } Z = \frac{\hat{\beta}_1 - \hat{\beta}_2 - (\beta_1 - \beta_2)}{SE(\hat{\beta}_1 - \hat{\beta}_2)} \sim N(0, 1)$$

where $SE(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_1 - \hat{\beta}_2)}$

### Nature of the Test

(i) For the alternative hypothesis $H$ ... $\beta - \beta_0$ will be accepted ... for the given sample ... and will be rejected otherwise.

(ii) For the alternative hypothesis $H$ ... be accepted if for the given sample ... when $t$ ...

(iii) For the alternative hypothesis $H$ ... be accepted if for the given sample ... (i.e. when $t$ ... in each case $t$ denotes the ...

### Confidence Interval of $(\beta_1 - \beta_2)$ :

At $100\,\alpha\%$ level of significance the confidence interval of $(\beta_1 - \beta_2)$ would be given by,

$$(\hat{\beta}_1 - \hat{\beta}_2) \pm t_{\alpha_2, n-3} SE(\hat{\beta}_1 - \hat{\beta}_2)$$

i.e.   $P[-t_{\alpha_2, n-3} < t \le t_{\alpha, n-3}] = 1 - \alpha$

or,   $P\left[ t_{\alpha_2, n-3} \le \dfrac{(\hat{\beta}_1 - \hat{\beta}_2) - (\beta_1 - \beta_2)}{SE(\hat{\beta}_1 - \hat{\beta}_2)} < t_{\alpha, n-3} \right] = 1 - \alpha$

or,   $P\left[ (\hat{\beta}_1 - \hat{\beta}_2) - t_{\alpha, n-3} SE(\hat{\beta}_1 - \hat{\beta}_2) \le (\beta_1 - \beta_2) \le \hat{\beta}_1 - \hat{\beta}_2 + t \ldots \right]$

$= 1 - \alpha$

where $(1 - \alpha)$ is the confidence coefficient.

### Case 5 - Confidence interval for $\sigma_u^2$:

Under the normality assumption, the variable

$$\chi^2 = \frac{RSS}{\sigma_u^2} = \frac{\Sigma e_t^2}{\sigma_u^2} = (n-3)\frac{\hat{\sigma}_u^2}{\sigma_u^2}$$

follows a $\chi^2$ (chi square) distribution with $df = n - 3$ [where $\hat{\sigma}_u^2 = \Sigma e_t^2 / (n-3)$, is an unbiased estimator of $\sigma_u^2$]

AN INTRODUCTION TO ...

Therefore we can use ... to establish a confidence interval for $\sigma$ ...

At the level ... significance the confidence limits for $\sigma$, we can by ...

$$ \sum \quad \text{and} \quad \sigma = \sum \quad \text{where } t \text{ values are taken from the table } \nu $$

$= (n - 3)$

$$ F \quad t \quad \quad 12 $$

$$ \beta_t \quad \quad = t $$

$$ F(t \quad \sigma_\pi \quad = \quad \alpha $$

where ... is the confidence coefficient.

**Example 3.8.** The following table contains observations on the quantity demanded of a certain commodity $Y$, its price $(X_1$, of $1)$ and consumer's income $X$, in $).

| Y | 10 | 95 | 80 | 70 | 70 | 65 | 97 | 90 | ... | 8 |
|---|---|---|---|---|---|---|---|---|---|---|
| X | 5 | 7 | 6 | 6 | 8 | 7 | 5 | 4 | 3 | 9 |
| t | 1000 | 600 | 300 | 500 | 700 | 400 | 300 | 100 | 300 | 600 |

Assume a linear regression equation of the form

$$ Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u, \quad i = 1, 2, 3 \dots \quad 10 $$

(i) Find the OLS estimators of $\beta_0$, $\beta_1$ and $\beta_2$ i.e. $\hat\beta_0$, $\hat\beta_1$ and $\hat\beta_2$

(ii) Find $R^2$ and adjusted $R^2(\bar{R})$

(iii) Find $var(\hat\beta_0)$, $var(\hat\beta_1)$ and $var(\hat\beta_2)$

(iv) Find $SE(\hat\beta_0)$, $SE(\hat\beta_1)$ and $SE(\hat\beta_2)$

(v) Write the regression results in the summary form.

(vi) Test $H_0: \beta_0 = 0$ against $H_1: \beta_0 \neq 0$ and find 95% and 99% confidence intervals for $\beta_0$

(vii) Test $H_0: \beta_1 = 0$ against $H_1: \beta_1 \neq 0$ and find 95% and 99% confidence intervals for $\beta_1$

(viii) Test $H_0: \beta_2 = 0$ against $H_1: \beta_2 \neq 0$ and find 95% and 99% confidence intervals for $\beta_2$

(ix) Test $H_0: \beta_1 = \beta_2$ against $H_1: \beta_1 \neq \beta_2$ and find 95% and 99% confidence intervals for $(\beta_1 - \beta_2)$

(x) Construct 95% and 99% confidence intervals of $\sigma_u^2$

| $n$ | $Y$ | $X_1$ | $X_2$ | $y_i = Y_i - \bar{Y}$ | $x_{1i} = X_{1i} - \bar{X}_1$ | $x_{2i} = X_{2i} - \bar{X}_2$ | $y_i^2$ | $x_{1i}^2$ | $x_{2i}^2$ | $x_{1i}y_i$ | $x_{2i}y_i$ | $x_{1i}x_{2i}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 100 | 5 | 1000 | 20 | 1 | 200 | 400 | 1 | 40,000 | 20 | 4000 | 200 |
| 2 | 75 | 7 | 600 | 5 | 1 | 200 | 25 | 1 | 40,000 | 5 | 1000 | -200 |
| 3 | 80 | 6 | 1200 | 0 | 0 | 400 | 0 | 0 | 160,000 | 0 | 0 | 0 |
| 4 | 70 | 6 | 500 | -10 | 0 | 300 | 100 | 0 | 90,000 | 0 | 3000 | 0 |
| 5 | 50 | 8 | 300 | 30 | 2 | 500 | 900 | 4 | 250,000 | -60 | -3000 | 1000 |
| 6 | 65 | 7 | 400 | -15 | 1 | -400 | 225 | 1 | 160,000 | -5 | 6000 | -400 |
| 7 | 90 | 5 | 1300 | 10 | -1 | 500 | 100 | 1 | 250,000 | -7 | 5000 | -500 |
| 8 | 100 | 4 | 1100 | 20 | 2 | 300 | 400 | 4 | 90,000 | -40 | 6000 | 600 |
| 9 | 110 | 3 | 1300 | 30 | 3 | 500 | 900 | 9 | 250,000 | 90 | 9000 | 900 |
| 10 | 60 | 9 | 300 | 20 | 3 | 500 | 400 | 9 | 250,000 | 60 | 9000 | 900 |
| $n=10$ | $\Sigma Y_i$ 800 | $\Sigma X_{1i}$ 60 | $\Sigma X_{2i}$ 8000 | $\Sigma y_i$ 0 | $\Sigma x_{1i}$ 0 | $\Sigma x_{2i}$ 0 | $\Sigma y_i^2$ 3450 | $\Sigma x_{1i}^2$ 30 | $\Sigma x_{2i}^2$ 880,000 | $\Sigma x_{1i}y_i$ 400 | $\Sigma x_{2i}y_i$ 68000 | $\Sigma x_{1i}x_{2i}$ 2900 |

$$\bar{Y} = \frac{\Sigma Y_i}{n} = \frac{800}{10} = 80, \qquad \bar{X}_1 = \frac{\Sigma X_{1i}}{n} = \frac{60}{10} = 6 \text{ and } \bar{X}_2 = \frac{\Sigma X_{2i}}{n} = \frac{8000}{10} = 800$$

**Solution**

The observations are in ...

$\hat{\beta}$ ...

$[\beta_1 \quad ...]$

Similarly $\hat{\beta}$ ...

$\beta \cdot ... = 3$

When $\beta_0$ and $\beta_1$ are known, $\beta_2$ can be obtained from the relation

$\beta_0 = \beta \bar{Y} - \beta \bar{X}$

... $\hat{\beta}_0$ ...

Thus we have $\hat{\beta}_0 = ...$ $\hat{\beta} = ...$ and $\hat{\beta}_2 = ...$

ii) We know that $R^2 = \frac{ESS}{TSS}$ ...

$R^2 = 0.894$. This means that price and income can mainly explain 89.4% variation in demand out of total variation of 100%.

Now, adjusted $R = \bar{R}^2 = ...$ Here $n = ...$

$\bar{R}^2 = (1 - 0.894) \cdot \frac{...}{...} ...$

$R^2 = 0.894$ and adjusted $R = \bar{R}^2 = 0.8637$

iii. We know that $var(\beta_1) = \frac{\sigma_u^2 \Sigma x_2^2}{\Sigma x_2^2 \Sigma x_2^2 - (\Sigma x_1 x_2)}$

Here $\sigma_u^2$ is unknown and is replaced by an unbiased estimator $\hat{\sigma}_u^2 = \frac{...}{...}$

We know that $R^2 = 1 - \frac{\sum u_i^2}{\sum y_i^2} = ...$

...

...

...

var(...) = 6.55

Again, var($\beta_1$) ... $\Sigma_1 \quad \Sigma c_{2i}^2 \quad \Sigma ...$

$$\frac{52.24 \times 10}{10 \times 580,000 - (-5900)} \quad ... = ... 24$$

var($\beta_2$) = 0.000124

Again, var($\beta_0$) = $\frac{\sigma_u^2}{n} + \bar{X}_1^2 \, \text{var}(\beta_1) - 2\bar{X}_1 \bar{X}_i \, \text{cov}(\beta_1 \beta_i) - ...$

We know that $\sigma_u^2 = 52.24$, $n = ...0$, $\bar{X}_1 = 6$, $\bar{X}_i = 800$, var(...) = 6.55, var(...) = 0.000124

and $\text{cov}(\hat{\beta}_1 \, \hat{\beta}_2) = \frac{-\hat{\sigma}_u^2 \Sigma c_{1i} x_{2i}}{\Sigma x_{1i}^2 \, \Sigma x_{2i}^2 - (\Sigma x_{1i} x_{2i})^2}$

$$= \frac{52.24 \times (-5900)}{30 \times .580,000 - (-5900)^2} = \frac{308216}{12590000} = 0.0245$$

cov($\hat{\beta}_1 \, \hat{\beta}_2$) 0.0245. We now put these values in the expression of $\beta_0$ and we get,

$$\text{var}(\hat{\beta}_0) = \frac{52.24}{10} + (6)^2 \times 6.55 + 2 \times 6 \times 800 \times 0.0245 - (800)^2 ... = ...$$

$$= 5.224 + 235.8 + 235.2 + 79.36 = 555.58$$

(iv) We know that $SE(\beta_1) = \sqrt{\text{var}(\beta_1)}$

$$SE(\hat{\beta}_1) = \sqrt{\text{var}(\hat{\beta}_1)} = \sqrt{6.55} = 2.5592$$

The regression results in summary form

$R^2 = 0.894$, Adjusted $R^2 = \bar{R}^2 = 0.8617$

We have to test the null hypothesis $H_0 : \beta_0 = 0$ against the alternative $H_1$. The appropriate test statistic under $H_0 : \beta_0 = 0$ would be

$$t = \frac{\hat{\beta}_0}{SE(\hat{\beta}_0)} \sim t_n$$

Here observed $t = \dfrac{\hat{\beta}_0}{SE(\hat{\beta}_0)} = \dfrac{111.70}{...} = 4.719$

The null hypothesis $H_0 : \beta_0 = 0$ will be accepted if ... for the given sample ... and will be rejected otherwise.

When $\alpha = ...$    $t_{0.05,10} = t_{0.025,7} = 2.365$

$t_{0.025,7} = 2.365$ (Table value)

Thus we see that $t$ observed $= 4.719$ does not lie in the interval $-2.365 ...$ and $+ 2.365 ...$ Hence $H_0 : \beta_0 = 0$ is rejected and $H_1 : \beta_0 \neq 0$ is accepted at 5% level of significance.

Similarly, when $\alpha = 0.01$ ... $t = t_{0.005} = 3.499$. Here we see that $t$ observed $= 4.719$ does not lie in the interval $-3.499$ and $3.499$. Hence $H_0 : \beta_0 = 0$ is rejected and $H_1 : \beta_0 \neq 0$ is accepted at 1% level of significance.

We know that $(1-\alpha)\%$ confidence interval of $\beta_0$ would be

$$P\left[ ... \leq ... \leq ... \right] = 1 - \alpha$$

or, $P\left[ ... \leq \dfrac{\hat{\beta}_0 - \beta_0}{SE(\hat{\beta}_0)} \leq t_{...} \right] = 1 - \alpha$

or, $P\left[ \hat{\beta}_0 - t_{...} SE(\hat{\beta}_0) \leq \beta_0 \leq \hat{\beta}_0 + t_{...} \right] = 1 - \alpha$

when $\alpha = 0.05$, $P[\hat{\beta}_0 - t_{0.025,7} \times 23.570 \leq \beta_0 \leq \hat{\beta}_0 + t_{0.025,7} \times 23.570] = 1 - 0.05 = 0.95$

or, $P[111.70 - 2.365 \times 23.570 \leq \beta_0 \leq 111.70 + 2.365 \times 23.570] = 0.95$

or, $P[55.957 \leq \beta_0 \leq 166.743] = 0.95$

95% confidence intervals of $\beta_0$ are 55.957 and 166.743

Similarly, when $\alpha = 0.01$, then $100(1-\alpha)\% = 99\%$

99% confidence intervals of $\beta_0$ would be

$$\hat{\beta}_0 \pm t_{\alpha/2, n-3} SE(\hat{\beta}_0)$$

or, $\hat{\beta}_0 - t_{0.025} SE(\hat{\beta}_0)$

.1 $70 \pm 3.499 \times 23.570$

or, $70 \pm 82.474$

or, $-10.9286$ and $194.1714$

So 99% confidence intervals of $\beta_0$ are $-10.9286$ and $194.1714$.

So, to test the null hypothesis $H_0: \beta_1 = 0$ against the alternative $H: \beta_1 \ne 0$ the test statistic under $H_0: \beta_1 = 0$ would be $t$ observed $= \dfrac{t_0}{SE(\hat{\beta}_1)}$

Here observed $= \dfrac{\hat{\beta}_1}{SE(\hat{\beta}_1)} = \dfrac{7.19}{2.5592} = 2.8094$

Now, $H_0: \beta_1 = 0$ will be accepted if for the given sample $|t \text{ observed}| \le t_{\alpha/2, n-3}$ and will be rejected otherwise.

When $\alpha = 0.05$ $t_{\alpha/2, n-3} = t_{0.025, (10-3)} = t_{0.025, 7} = 2.365$

and when $\alpha = 0.01$, $t_{\alpha/2, n-3} = t_{0.025, 7} = 3.499$

Here we see that $t$ observed $= 2.8094$ does not lie in the interval $-2.365$ and $2.365$ and hence $H_0: \beta_1 = 0$ is rejected at 5% level of significance. But $t$ observed $= 2.8094$ lies in the interval $-3.499$ and $3.499$ and hence $H_0: \beta_1 = 0$ is accepted at 1% level of significance.

(R) 95% confidence limits to $\beta_1$ would be

$$\hat{\beta}_1 \pm t_{\alpha/2, n-3} SE(\hat{\beta}_1)$$

when $\alpha = 0.05$, then 95% confidence limits to $\beta_1$ would be

$$\hat{\beta}_1 \pm t_{0.025, 7} SE(\hat{\beta}_1)$$

or, $7.19 \pm 2.365 \times 2.5592$

or, $7.19 \pm 6.0525$ or, $13.2425$ and $1.135$

So, 95% confidence limits to $\beta_1$ are $13.2425$ and $1.135$.

Similarly, when $\alpha = 0.01$ then 99% confidence limits to $\beta_1$ would be

$$\hat{\beta}_1 \pm t_{0.025, 7} SE(\hat{\beta}_1)$$

or, $7.19 \pm 3.499 \times 2.5592$

or, $7.19 \pm 8.9546$ or, $16.1446$ and $1.7646$

So, 99% confidence limits to $\beta_1$ are $16.1446$ and $1.7646$.

(viii) To test the null hypothesis $H_0: \beta_2 = 0$ against the alternative $H: \beta_2 \ne 0$ the appropriate test statistic under $H_0: \beta_2 = 0$ would be

$$t = \frac{\hat{\beta}_2}{SE(\hat{\beta}_2)} \sim t_{n-3}$$

... and will be rejected otherwise.

When ...

and when ...

Here we see that observed $\dfrac{\beta_1}{\text{...}}$ ...

... 

... confidence intervals of $\beta$ would be

$$\hat{\beta} \pm \text{...} SE(\hat{\beta})$$

when $\alpha$ ... the 95% confidence intervals of $\beta$ would be

$$\hat{\beta} \pm \text{...}$$

or, $0.0143 \pm 1.385 \times 0.01$ (

or ... or $0.00\text{...}$ and $0.0409$

So 95% confidence intervals of $\beta$ are ... and ...

When $\alpha$ ... then 99% ... 99% confidence intervals of $\beta$ would be

$$\hat{\beta}_2 \pm t_{0.005,7}\, SE(\hat{\beta}_2)$$

or ... $499$ ...

or ... $0.085$ or ... $45$ and $0.051$

So 99% confidence intervals of $\beta_2$ are ... $0.0245$ and $0.051$

(ix) To test the null hypothesis $H_0$: $\beta_1 = \beta_2$ against the alternative $\beta_1 \neq \beta_2$. The appropriate test statistics under $H_0$: $\beta_1 = \beta_2$ would be

$$t = \frac{\hat{\beta}_1 - \hat{\beta}_2}{SE(\hat{\beta}_1 - \hat{\beta}_2)} \sim t_{n-2}$$

Now $t$ (observed) $= \dfrac{\hat{\beta}_1 - \hat{\beta}_2}{SE(\hat{\beta}_1 - \hat{\beta}_2)}$

Since $\beta_1 = 7.0$, $\beta_2 = 0.0143$, $\text{var}\beta_1 = 6.55$, $\text{var}\beta_2 = 0.000124$ and

$SE(\hat{\beta}_1 - \hat{\beta}_2) = \sqrt{\text{var}(\hat{\beta}_1) + \text{var}(\hat{\beta}_2) - 2\text{cov}(\hat{\beta}_1, \hat{\beta}_2)}$

$\qquad = \sqrt{6.55 + 0.000124 - 2 \times 0.0245}$

$\qquad = \sqrt{6.55 + 0.000124 - 0.049} = \sqrt{6.50\text{...}} = 2.549$

Now $n$ _____ will be accepted _ for the _____

will be rejected otherwise.

When $\alpha = 0.05$, $t_{\alpha/2, n-3} = t_{0.025, 7} = 2.365$

and when $\alpha = $ ____ $t_{\alpha}$ __ ____

Here we see that $t$ observed, ____ $2.365$ does not lie _____

and hence $H_0$: $\beta_1$ _ $\beta_2$ is rejected and if ____ $\beta_1 \neq \beta_2$ is accepted at 5% _____

significance.

Will we see that $t$ observed _____ $2.k2n3$ lies in the interval _____

the null hypothesis $H_0$: $\beta_1$ _ $\beta_2$ is accepted at ____ _____

Now $100(1 - \alpha)$% confidence intervals of $(\beta_1 - \beta_2)$ would be

$$\beta_1 - \beta_2) \pm t_{\alpha/2, n-3} \cdot SE(\hat{\beta}_1 - \hat{\beta}_2)$$

when $\alpha = 0.05$, $100(1 - \alpha)$% = 95% confidence intervals for ____ $\beta$ would be

$$\hat{\beta}_1 - \hat{\beta}_2 \pm t_{0.025, 7} \cdot SE(\hat{\beta}_1 - \hat{\beta}_2)$$

or, $-7.9 - 0.6343 - 2.365 \times 2.549$

or, $-7.3043 + 6.0283$ or, $13.2326$ and __ __

So, 95% confidence intervals of $(\beta_1 - \beta_2)$ are $-13.2326$ and $1.376$

When $\alpha = 0.0$ then $100(1 - \alpha)$% = 99% confidence intervals of $\beta$ __ $\beta$ would

be

$$\hat{\beta}_1 - \hat{\beta}_2 \pm t_{\alpha/2, n-3} \cdot SE(\hat{\beta}_1 - \hat{\beta}_2)$$

or, $-7.9 - 0.0 - 3) \pm t_{0.005} - + 2.549$

or, $-7.2043 \pm 3.499 \times 2.549$

or, $7.2043 \pm 8.9189$ or, $16.232$ and $1.7146$

So, 99% confidence intervals of $(\beta_1 - \beta_2)$ are $-16.1232$ and $1.7146$

(b) We have to construct 95% and 99% confidence intervals of $\sigma_u^2$

We know that $100(1 - \alpha)$% confidence intervals of $\sigma_u^2$ would be given by

$$(n - 3) \frac{\sigma_u^2}{\chi^2_{\alpha/2, n-3}} \quad \text{and} \quad (n - 3) \frac{\sigma_u^2}{\chi^2_{1 - \alpha/2, n-3}} \quad \text{where } \chi^2 \text{ values are taken from the table}$$

with d.f $= n - 3$

When $\alpha = 0.05$ $\chi^2_{\alpha/2, n-3} = \chi^2_{0.025, 7} = 16.013$

and $\chi^2_{1-\alpha/2, n-3} = \chi^2_{0.975, 7} = 1.690$

and $\qquad = \dfrac{52.24}{\;} = \dfrac{7 \times 52.24}{\;} = 216.17$

and $\Sigma \qquad$

### 3.10 Analysis of Variance (ANOVA) in a Multiple Linear (Three Variable) Regression Model

Yet another item that is often presented in connection with the three variable linear regression model is the analysis of variance, i.e. the break down of the total sum of squares (TSS) into explained sum of squares (ESS) and the residual sum of squares (RSS).

The estimated three variable linear regression line, where the regression equation is $\hat{Y} = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$, is given by $\hat{Y} = \beta_0 + \beta_1 X_1$ and

$$\bar{Y} = \beta_0 + \beta_1 \bar{X}_1 + \beta_2 \bar{X}_2$$

Taking deviations from mean we have

$$\hat{Y}_i - \bar{Y} = \beta_0 + \beta_1 X_{1i} - \beta_0 - \beta_1 \bar{X}_1 + \beta_2 X_{2i} - \beta_2 \bar{X}_2$$

or $\quad \hat{y}_i = \beta_1 (X_{1i} - \bar{X}_1) + \beta_2 (X_{2i} - \bar{X}_2)$

or $\quad \hat{y}_i = \beta_1 x_{1i} + \beta_2 x_{2i}$, where $x_{1i} = X_{1i} - \bar{X}_1$ and $x_{2i} = X_{2i} - \bar{X}_2$

Now, error of estimate $e_i = Y_i - \hat{Y}_i$

Now, $\Sigma y_i^2 = \Sigma (\hat{y}_i + e_i)^2 = \Sigma \hat{y}_i^2 + 2\Sigma \hat{y}_i e_i + \Sigma e_i^2$

$\Sigma y_i^2 = \Sigma \hat{y}_i^2 + \Sigma e_i^2$ as $\Sigma \hat{y}_i e_i = 0$ by assumption

i.e. TSS = ESS + RSS

with $K = 3$ as there are three ...

## ANOVA TABLE
### Table 3.4

| Source of variation | Sum of squares | Degrees of freedom | Mean sum of squares | | F | | |
|---|---|---|---|---|---|---|---|
| Explained | $ESS$ $\Sigma\hat{y}^2$ | | | | | | |
| (between) | $\beta \Sigma x_{2i} y$ | $k$ | $MSE$ | | $\dfrac{M.E}{M.R}$ | | |
| | $\beta \Sigma x_{3i} y$ | | | | with d.f | | |
| | | | | | $k, n-k$ | | |
| Residual (within) | $RSS$ $\Sigma e_i^2$ | $n-k$ | $\dfrac{RSS}{n-k}$ $MSR$ | | | | |
| Total | $TSS$ $\Sigma y^2$ | $n-1$ | | | | | |

In case of a three variable linear regression model, there are three parameters and hence $k = 3$.

The test aims at finding out whether the explanatory variables $X_2$ and $X_3$ do actually have any significant influence on the dependent variable $Y$. Forming the test of the overall significance of the regression implies testing the null hypothesis $H_0$: $\beta_2 = \beta_3 = 0$ against the alternative hypothesis $H_1$: not all $\beta_i$ are zero. We may use the test statistic

$$F = \frac{MSE}{MSR} = \frac{\Sigma \hat{e}^2 \ / \ (k-1)}{\Sigma e_i^2 \ / \ (n-k)}$$

$$= \frac{\Sigma \hat{y}^2 \ / \ (k-1)}{\Sigma e_i^2 \ / \ (n-k)} \quad \text{with d.f} = k-1,\ n-3 \text{ (Here } K = 3)$$

Now we have to compare $F^*$ with the table value of $F$ with d.f $= 2, n-3$. If it is found that $F^* > F_{\alpha,\ 2,\ n-3}$ (Table value) we reject the null hypothesis at $100\ \alpha\%$ level of significance ($\alpha = 0.01$ or $0.05$ usually), i.e. we accept that the regression is significant and not all $\beta$'s are zero.

$ F^* ...$ we accept the null hypothesis that ... ... ... ... regression is not significant

**Note  Relation between $R$ and $F$**

There is an intimate relationship between the coefficient ... multiple ... $R^2$ and the $F$ test used in the analysis ... variance ... ... normal ... ... for the disturbances $u$ and the null hypothesis ... ... ... we get $H_0$

$$F^* = \frac{ESS ...}{RSS ...}$$

... distributed as the $F$ distribution with d.f. ... and $n$ ... where $K$ ... is here if three parameters ... using the constant intercept term ... in a three variable ... regression model

Now we can write $F = F^* = \dfrac{ESS \ (K-1)}{RSS \ ... \ A}$

... ... ...

$$F^* = \frac{R}{A \ ... \ n \ A} = R$$

It should be noted that here $K$ = number of parameters in the ... regression ... ... $K$ = ... when there are two explanatory variables When $R^2 = 1, F^* = ...$ the larger the $R^2$ the greater the $F^*$ value

... the time when $R^2 = ...$ is ... Thus the $F$ test which is a measure of ... when ... significance of the estimated regression is also a test of significance ... $R^2$ It other words, testing the null hypothesis $H_0: \beta_2 = \beta_3 = 0$ is equivalent to test the null hypothesis that population $R^2$ is zero. The ANOVA table can also be written expressed in terms of $R^2$ as shown below

**ANOVA TABLE in terms of $R^2$**
**Table 3.5**

| Source of variation | Sum of squares SS | Degrees of freedom d.f | Mean sum of squares MS | F | |
|---|---|---|---|---|---|
| | | | | Observed | Tabulated |
| Explained (between) | $ESS = \Sigma \hat{y}^2$ $= R^2 \Sigma y^2$ | $K-1$ | $ESS/(K-1)$ $= MSE$ | $F = \dfrac{MSE}{MSR}$ with d.f | |
| Residual (within) | $RSS = \Sigma e^2 =$ $= (1-R^2) \Sigma y^2$ | $n-K$ | $RSS/(n-K)$ $= MSR$ | $= (K-1, n-K)$ | |
| Total | $TSS = \Sigma y^2$ | $n-1$ | | | |

Example 3.9 ... $H_1$, $\beta_1$ and $\beta_2$ are not zero.

Solution ...

... $\dfrac{\beta_2}{SE \, \beta_2}$ and $t = \dfrac{\beta}{SE \, \beta}$ ...

ANOVA Table

**ANOVA TABLE relating to demand for a commodity**
**Table 3.6**

| Source of variation | Sum of squares (SS) | Degrees of freedom (D) | Mean sum of squares (MS) | | |
|---|---|---|---|---|---|
| regression | ESS 2 | K — 2 | MSE $\dfrac{E}{K}$ | $F$ $F^*$ | |
| between ... | 3086.5 | | 543.35 | $\dfrac{M}{M \, R}$ | |
| Residual | RSS = $\Sigma e_i^2$ | $n - K$ 7 | MSR $\dfrac{RSS}{n-K}$ | with d.f. | |
| (within) | 363.5 | | 51.12 | $(K-1, n-K)$ 2, 7 | |
| Total | $\Sigma y^2 = 3450$ | $n - 1 = 9$ | | | |

Here the sample size, $n = 10$. number of parameters $K = 2$

$K - 1 = 2$ and $n - K = 10 - 3 = 7$

From Example 3.8 we have obtained the results

$ESS = \Sigma \hat{y}_i^2 = \beta_1 \Sigma x_{1i} y + \beta_2 \Sigma x_{2i} y_i = 3086.5$

$RSS = \Sigma e^2 = 363.5$ and $TSS = \Sigma y_i^2 = 3450$

Now the null hypothesis $H_0 : \beta_1 = \beta_2 = 0$ will be rejected if for the given sample

$F = F^* \text{ (observed)} = \dfrac{MSE}{MSR}$ [with d.f $(K - 1) = 2$ and $(n - K) = 7$]

is greater than the table value of $F$ with d.f. $(K - 1) = 2$ and $n - K = 7$. From the table value we see that $F_{0.05, \, 2,7} = 7.74$ and $F_{0.01, \, 2,7} = 9.55$

Here we see that $F \text{ (observed)} = F^* = 29.72$ and $F_{.05, 2,7} = 7.74$

$F^* > F_{0.5, \, 2,7}$

So, at 5% level of significance the null hypothesis $H_0 : \beta_1 = \beta_2 = 0$ will be rejected for the given sample.

We also see that $F^* = 29.72 > F_{0.01, \, 2,7} = 9.55$. This means that at 1% level of significance the null hypothesis $H_0 : \beta_1 = \beta_2 = 0$ will be rejected for the given sample.

Thus both at 1% and 5% levels of significance we may claim that the coefficients of the regression equation are not zero.

It should be noted that we can also construct the ANOVA table in the following Example 3.8 we are showing the ANOVA table below in terms of $R^2$.

### ANOVA TABLE in terms of $R^2$
### Table 3.7

| Source of variation | Sum of squares SS | Degrees of freedom (d.f.) | Mean sum of squares (MS) | F Observed | F tabulated |
|---|---|---|---|---|---|
| Explained (between) | $\Sigma SS \cdot R^2$ $3094 \cdot WSL$ $3044 \cdot 30$ | $k-1=2$ | $\frac{ESS}{(k-1)}$ $= MSE$ $\frac{3094.30}{2}$ $542.5$ | $F^*$ $= \frac{MSE}{MSR}$ $= \frac{1542.5}{52.242}$ $29.52$ | $F^*$ 7.74 |
| Residual (within) | $RSS$ $\Sigma (1-R^2)$ $\cdot 3450 \cdot (1-R^2)$ $365.70$ | $n-k=7$ | $\frac{RSS}{(n-k)}$ $MSR$ $\frac{365.70}{7}$ $52.242$ | 29.52 | $F^*$ $0.55$ |
| Total | $TSS$ $\Sigma y^2 = 3450$ | $n-1=9$ | — | | |

From Example 3.8 we have seen that $n = 10$, $k = 3$, $\Sigma y^2 = 3450$ and $R^2 = 0.894$. Here also we see that $F = F^* = 29.52 > F_{0.05} = 7.74$ and $F^* = 29.52 > F_{0.01} = 9.55$.

Thus the null hypothesis $H_0 : \beta_1 = \beta_2 = 0$ is rejected both at 1% and 5% levels of significance.

## 3.8. The Cobb-Douglas Production Function : More on Functional Form

In Section 2.18 we showed how with appropriate transformations we can convert non-linear relationships into linear ones so that we can work within the framework of classical linear regression model. We consider the Cobb-Douglas Production function which shows a three variable non-linear relation. The Cobb-Douglas Production function, in its stochastic form, may be expressed as

$$Y = \beta_0 X_1^{\beta_1} X_2^{\beta_2} U_i$$

where $Y$ = output, $X_1$ = labour input, $X_2$ = capital input, $U$ = Stochastic disturbance term, $\beta_0$ = constant technological parameter.

Taking log on both sides of the ...

$\log Y = \log \beta_0 + ... $

... $\beta_1 ...$

...

We ... now apply the ... ... ... thing the OLS estimators of the Cobb-Douglas production function

... is the partial ... ... with respect of labour input i.e.

$$\log Y$$
$$\log X_1$$

(ii) Likewise ... is the partial ... ... ... ... input

$$\frac{\log Y}{\log X_2} = \beta_2$$

(iii) The sum $(\beta_1 + \beta_2)$ gives information about the ... ... displays IRS, CRS and DRS according as $\beta_1 + \beta_2 \lesseqgtr 1$.

**Example 3.10.** A production function is specified as $Y_i = \beta_0 X_{1i}^{\beta_1} X_{2i}^{\beta_2} U_i$, where $i = 1, 2, \dots, n$

$Y$ = output, $X_1$ = labour input, $X_2$ = capital input, $U$ = Stochastic disturbance term, $n$ = sample size. The corresponding Log-linear form of the production function is given as

$$\log Y_i = \log \beta_0 + \beta_1 \log X_{1i} + \beta_2 \log X_{2i} + \log U_i$$

or $y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + u_i$, $u_i \sim N(0, \sigma_u^2)$

On the basis of a sample size of 23 the following results are given : $\hat{\alpha} = 4.0$, $\hat{\beta}_1 = 0.7$, $\hat{\beta}_2 = 0.2$, $RSS = 1.4$, $TSS = 10$, $var(u) = 0.6084$, $var(\hat{\beta}_1) = var(\hat{\beta}_2) = 1.35$

(i) Write the estimated regression equation.

(ii) Find the value of $R^2$

(iii) Find $SE(\hat{\alpha})$, $SE(\hat{\beta}_1)$, and $SE(\hat{\beta}_2)$

(iv) Find $\hat{\sigma}_u^2$

(v) Find the 95% confidence intervals for $\alpha$, $\beta_1$, $\beta_2$ and $\sigma_u^2$

(vi) Test the hypothesis $\beta_1 = 1.0$ and $\beta_2 = 0$ separately at the 5% significance level

**Solution :** (i) The estimated regression equation can be written as

$$\hat{y}_i = \hat{\alpha} + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i}$$

or, $\hat{y}_i = 4.0 + 0.7 x_{1i} + 0.2 x_{2i}$

Y-60-11

(i) The value of multiple coefficient of determination, given by $R^2$ ...

where $ESS + RSS = $ ... $R^2 = ESS/$ ...

We know that $SE(\hat{\beta}_1) = $ ... varies ... ... = 0 ...

similarly, if $\beta_1$ ... varies ... and $\sigma = $ ... 0.10?

like $SE_{\hat{\beta}} = $ ... and ... $SE(\hat{\beta}) = 0.10$?

(iv) We have to find out the value of OLS estimator of $\sigma_u = E$ ... $\hat{u}$

We know that $\sigma_u = \dfrac{\sum \hat{u}^2}{n - 3} = \dfrac{RSS}{n - 3}$

Here $RSS = $ ... 4 and $n = 23$

$$\sigma_u = \frac{RSS}{n-3} = \frac{4}{20} = 0.0 ...$$

(v) We have to find out the 95% confidence intervals for $\alpha$, $\beta_1$, $\beta_2$ and $\sigma_u$. Using the $t$ distribution with d.f. $(n - 3) = (24 - 3) = 20$, we can get the confidence intervals for $\alpha$, $\beta_1$ and $\beta_2$ as

**For $\alpha$:** $\hat{\alpha} \pm t_{0.025, 20} SE(\hat{\alpha}) = 4.0 \pm 2.086 \times ... = ...$

... $+ 37 + 63 = 1 ...$

... $t_{0.025,20} = 2.086$ as $n - 3 = 23$ ...

**For $\beta_1$:** $\hat{\beta}_1 \pm t_{0.025} SE(\hat{\beta}_1) = 0.7 \pm .086 \times 0.10$

$$= 0.7 \pm 0.21 = (0.49, 0.91)$$

and **for $\beta_2$:** $\hat{\beta}_2 \pm t_{0.025} SE(\hat{\beta}_2) = 0.2 \pm 2.086 \times 0.02$

$$= 0.2 \pm 0.21 = (-0.01, 0.41)$$

Again, 95% confidence intervals for $\sigma_u$ would be

$$P\left[ (n-3)\frac{\hat{\sigma}_u^2}{\chi^2_{\alpha/2, n-3}} \leq \sigma_u^2 \leq (n-3)\frac{\hat{\sigma}_u^2}{\chi^2_{1-\alpha/2, n-3}} \right] = 1 - \alpha$$

When $\alpha = 0.05$, $n = 23$, $\hat{\sigma}_u^2 = 0.07$

$$P\left[ 20 \cdot \frac{0.07}{\chi^2_{0.025, 20}} \leq \sigma_u^2 \leq 20 \cdot \frac{0.07}{\chi^2_{0.975, 20}} \right] = 1 - 0.05 = 0.95$$

or $\quad P\left[ \dfrac{1.4}{34.70} \leq \sigma_u^2 \leq \dfrac{1.4}{9.59} \right] = 0.95$

or $\quad P\left[ 0.04 \leq \sigma_u^2 \leq 0.146 \right] = 0.95$

95% confidence intervals for $\sigma_u$ are 0.04 and 0.146.

to test the null hypothesis $H_0 : \beta_1$ ... the alternative $H_a : \beta_1 \neq 1$, the

... test statistic under the $H_0 : \beta$ ...

$$t \text{ (observed)} = \frac{\hat{\beta}_1 - 1}{SE(\hat{\beta}_1)} \sim t_{...}$$

$$\therefore t \text{ (observed)} = \frac{\hat{\beta}_1 - 1}{SE(\hat{\beta}_1)} = \frac{0.7 - 1}{0.102} \cdot$$

$H_0 : \beta_1$ ... will be ...

... and will be rejected otherwise

When $\alpha = 0.05$, $t_{\alpha/2, n-2} = t_{0.025, 20} = 2.086$

... we see that $t$ (observed) $= -2.941$

... does not lie in the interval ... $\beta_1$ ... and ...

null hypothesis $H_0 : \beta_1 = 1.0$ is rejected at 5% level of significance

Again to test the null hypothesis $H_0 : \beta_2$ ... against the ... $H_a$ ...

the appropriate test statistic under $H_0 : \beta_2$ ... would be

$$t \text{ (observed)} = \frac{\hat{\beta}_2 - 0}{SE(\hat{\beta}_2)} \sim t_{n-3}$$

Here $t$ (observed) $= \frac{0.2}{0.102} = 1.960$

Now, $H_0 : \beta_2 = 0$ will be accepted if for the given sample $t$ ... observed ...

... and will be rejected otherwise

When $\alpha = 0.05$, $t_{\alpha/2, n-3} = t_{0.025, 20} = 2.086$

Here we see that $t$ (observed) $= 1.960$ lies in the interval $-2.086$ and $2.086$ and hence
$H_0 : \beta_2 = 0$ is accepted at 5% level of significance.

## 3.12. Prediction / Forecasting in the Multiple (Three-Variable) Regression Model

The formulas for prediction in multiple regression are similar to those in the case
of simple (two variable linear regression) regression, except that to compute the standard
error of the predicted value we need the variances and covariances of a ... regression
parameters. Here we will present the expression for the standard error in the case of
two explanatory variables.

Let the estimated regression equation be,

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2$$

Now consider the prediction of the value $Y_0$ of $Y$ given values $X_{10}$ of $X_1$ and $X_{20}$
of $X_2$, respectively. These could be values at some future date

Then we have $Y_0 = \beta_0 + \beta_1 X_{10} + \beta_2 X_{20} + u_0$

and $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 X_{10} + \hat{\beta}_2 X_{20}$

The prediction of ...

Thus the prediction ...

The variance of the prediction ...

where ...

When ... is first estimated, it is replaced by its estimated variance ... Thus ... is a confidence interval ... the prediction will be

$\hat{Y}_0 = t_{\alpha/2} \cdot SE \cdot \hat{Y}_0$ where $SE = ...$

**Example 1.11** Following Example ...

using the estimated regression ... Find the prediction of $Y$ for $X_1 = ...$ and $X_2 = 40$.

Using the results of ... $R$ and ... derived in Example 3.5, estimate the variance of the prediction error and the standard error of the prediction error.

(iii) Find 95% confidence interval for the prediction.

**Solution**

(i) The prediction of $Y$ can be obtained from the estimated regression equation.

$\hat{Y} = 11.70 + 9 X_0 + ...$

Here we put $X_{10} = 10$ and $X_{20} = 40$ and we get

$\hat{Y}_0 = 11.70 + 9 \times 10 + ... = 55.54$

The predicted value of $Y$ is $\hat{Y}_0 = 55.54$ when $X_1 = 10$ and $X_2 = 40$.

(ii) From example 3.5 we obtain ... var($\hat{\beta}$) = 0.3003 ... and $\hat{\beta}$ ...

$$\text{var}(\hat{Y}_0) = 32.24\left[1 + \frac{1}{10}\right] + (10 - 6)^2 + 6.11$$

$+ (10 - 6)(1400 - 1000) + 0.0 \, 45 + (1400 - 1000)^2 + 0.00034$

Here $X_{10} = 10$ and $X_{20} = 1400$, which are given...

$= 35.4 + 1...$

$= 35.464$...

$= 174.504$

$\text{var}(\hat{Y}_0) = 174.504$

and $SE(\hat{Y}_0) = \sqrt{\text{var} \ldots}$ ...

Variance of the prediction is ... of ... and standard error of prediction is ... 
.804

(b) We know that $(1 - \alpha)100\%$ confidence interval for prediction would be

$\hat{Y}_0 \pm t_{\alpha/2, n-1} SE(\hat{Y}_0)$ when $\alpha = 0.05$ ...

95% confidence interval for prediction would be

$\hat{Y}_0 \pm t_{0.025, 7} SE(\hat{Y}_0)$

or $58.56 \pm 2.365 \times 18.014$ or $58.56 \pm 42.60$ or $(15.96, \ldots)$ ...

95% confidence interval for prediction would be ... 5.96 and ...

## 3.13. Regression Analysis in Presence of Qualitative (Dummy) Variables

### 3.13.1. Meaning

In Section 1.10 we mentioned four types of variables that one generally encounters in empirical analysis. These are – ratio scale, interval scale, ordinal scale and nominal scale. The types of variables used in earlier sections were essentially ratio scale. In many cases we deal with models that may involve not only ratio scale variables but also nominal scale variables. Such variables are known as indicator variables, categorical variables, qualitative variables, or, dummy variables (Binary variables).

### 3.13.2. Nature of Dummy Variables

In regression analysis the dependent variable or regressand is frequently influenced not only by ratio scale variables (e.g. income, output, prices, costs, height, weight, temperature, etc.) but also by variables that are essentially qualitative or nominal scale, in nature, such as sex, race, colour, religion, nationality, geographical region, political upheavals and party affiliation. For example, holding all other factors constant, female workers are found to earn less than their male counterparts or non-white workers are

[Top portion of page too faded to read reliably]

Such models are also called **Analysis of variance (ANOVA) models**.

### Use of Dummy Variables

The dummy variables are used differently in the regression analysis. The dummy variables can be used for several purposes. We are explaining below some of the purposes for which a particular dummy variable is used.

#### (i) Dummy variables as proxies to Qualitative (Categorical) Factors

Dummy variables are sometimes used as proxies for qualitative factors in a regression regression. For example we consider a sample of a number of persons selected at random from the country, rural and urban and we want to estimate the demand for a particular item such as a bank out of income. It is known that town dwellers have higher demands than rural districts. Thus region is an important explanatory factor in this case. We may represent this factor by a dummy variable to which we may arbitrarily assign the value 1 for a town dweller and 0 for a person living in a rural area. The demand function can be written as

$$D = \beta_0 + \beta_1 Y + \beta_2 X + u$$

where Y = income and X = dummy variable for region. We put X = 1 for a town dweller ($\beta_0$) and X = 0 for a person living in rural area.

#### (ii) Dummy variables as proxies of Numerical Factors

Dummy variables may be used as proxies for quantitative factors when no observations on these factors are available or when it is convenient to do so. In example suppose we want to measure saving function $S = a + bY$ from a cross-section sample of consumers. Obviously, age is an important explanatory factor of the consumption and savings patterns of a consumer, since people become more thrifty with age. Although age is a quantitative factor we may approximate it by a dummy variable. We may divide the selected consumers in two age groups.

Group 1  People of 30–40 years of age

Group 2  People of 40 years and over

On the assumption that people become more thrifty as they grow old, the dummy variable for age may be assigned the value 0, if the person belongs to group 1 and the value 1 if the person belongs to group-2.

The saving function can be written in the form

$$S_i = \beta_0 + \beta_1 Y_i + \ldots$$

where $S$ = saving, $Y$ = ...

... dummy variable ...

where $D_i = 0$

**(iii) Dummy variables are used for measuring the shift of a function over time**

A shift of a function implies that the constant term changes while the other coefficients remain constant. Such shifts may be captured by the introduction of a dummy variable in the function.

For example suppose that we have data on the consumption for an economy for the period 1910-1968. During this period the economy faces two World Wars (1914-1918 and 1939-1945) and a deep depression (1929-1933). The abnormal conditions prevailing in these years have caused a shift of the consumption function due to war, depression, rationing, various controls and other factors. To capture this shift we may use a dummy variable say $Z$ which would assume the value 0 during the above 'abnormal' years and 1 in the other normal years. The consumption function takes the form

$$C_i = \beta_0 + \beta_1 Y_i + \beta_2 Z_i + u_i, (\beta_1 > 0)$$

where $C$ = consumption, $Y$ = income,

$Z$ = Dummy variable for the shift of the function

For a normal year the estimated form of the consumption function would be

$$\hat{C} = \hat{\beta}_0 + \hat{\beta}_1 Y + \hat{\beta}_2 = (\hat{\beta}_0 + \hat{\beta}_2) + \hat{\beta}_1 Y$$

and for an abnormal period it would be

$$\hat{C} = \hat{\beta}_0 + \hat{\beta}_1 Y$$

If we plot these two functions we can clearly see the shift in the consumption function during the abnormal (War and depression) years



**Fig. 3.1**

The slope of the consumption function (Fig 3.1) i.e. MPC is assumed to be the same both in normal and abnormal periods and hence the two regression lines are parallel (only intercept changes, slope remaining the same).

**(iv) Dummy variables are used for measuring the change of parameters (slope) over time**

It is usual that in certain period of time or in abnormal war depression years not only the income or the price rate cpi changes but also it is possible as well be expected to change. Elasticities and propensities change this change in the parameters. A variation can be captured by introducing appropriate dummy variables in the function.

We can write here the consumption function in the form

$$C_i = \beta_0 + \beta_1 Y_i + \beta_2 Z_{1i} + \beta_3 Z_{2i} + u_i$$

where    consumption, $Y$ income

$Z_1$ = Dummy variable    $\begin{cases} \text{abnormal years} \\ \text{normal years} \end{cases}$

$Z_2 \cdot Y_i$   Dummy variable   $\begin{cases} 0 \text{ for abnormal years (when } Z_1 = 0 \\ \text{for normal years (when } Z = \end{cases}$

Consequently for a normal period the estimated consumption function will be given by

$$\hat{C} = \beta_0 + \beta_1 Y$$   where for an abnormal year the estimated function would be

$$\hat{C} = \beta_1 Y$$

In this case both slope and intercept of the function will change

**(v) Dummy variables are used as proxies for the dependant variable**

In some cases the dependent variable of a function may be a dummy variable. For example suppose we want to measure the determinants of car ownership from a cross-section sample. Some people will have cars while others will not. Suppose that the determinants of the ownership depend on income and profession.

The functional relation can be written in the form

$$Y_i = \beta_0 + \beta_1 Y_i + \beta_2 Z_i + u_i$$

where $C$   car owners or non-owners

$= \begin{cases} \text{for car owners} \\ 0 \text{ for non owner} \end{cases}$

Here $C$ is taken as a dummy variable
$Y$   income
$Z$   a dummy variable for profession
$= \begin{cases} 1 \text{ if employed formally} \\ 0 \text{ if employed informally} \end{cases}$

It should be noted that if the dependent variable of a function is taken as a dummy variable, the disturbance term will be heteroscedastic and instead of OLS will not be appropriate there

**(vi) Dummy variables are used for seasonal adjustments of time series**

One of the most common use of dummy variables is in removing seasonal variations in time series. For example if we have quarterly data on retail sales we should adjust

... work out bases at Christmas ...
... when the sales on ...
... where the sales on the expense ...
... in my proper is the expense ...

... the quarters regions ...

... for the ...

$(3.$   $0 ...$

$w$ ...    $1 ...$

$\quad$ 3 in the third quarter

$D_{4}$   0 in all other quarters

## Dummy variable trap

It should be noted that we can ...
values for the fourth quarter and ...
of the terms of sums of squares and sum of ...
including the quarterly dummies) would be zero ...
which is introduced with values equal to 1 in all ...
constant intercept $\beta_0$.

If we apply OLS to the above quarter's model, the ...
will give seasonal effect for each of the three ...
are zero and the seasonal effect for the fourth quarter ...
$\beta_0$

In fact, when we introduce a large number of dummy variables ... the model, we
cannot obtain the OLS estimators of the parameters ... $(X'X)$ matrix may
be singular and $(X'X)^{-1}$ may not exist. This problem ... Dummy variable trap.

## Some illustrative Examples .

**Example 3.12.** Consider the following model showing consumption expenditure by
geographical region

$$Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + u_i$$

where   $Y_i$ = Average consumption expenditure ($) per person per 30 days in
State $i$

$$D_{1i} = \begin{cases} = 1, \text{ if the State is in the Eastern region of India} \\ = 0, \text{ otherwise (i.e., in other region of the country} \end{cases}$$

$$D_{2i} = \begin{cases} 1, \text{ if the State is in the North West region of the country} \\ 0, \text{ otherwise (i.e., in other region of the country)} \end{cases}$$

Using data for 17 States of India in 2006-07 the following results are obtained by
OLS method

$$\hat{Y}_i = 1097.38 \quad 241.04 D_{1i} \quad 30.09 D_{2i}$$

$$SE \quad (103.31) \quad (133.37) \quad (129.50)$$

$$t \quad (10.62) \quad (1.81) \quad (0.23)$$

The regression results show the mean percapita consumption expenditure ... in the eastern region the ... consumption about ₹ ... and that in the North-West central region ... is lower about ...

**Example 3.13.** We consider a model to show the different kinds of ...
literacy rates of States of India 2000-01

The model takes the form

$$Y_i = \beta_0 + \beta_1 D_{1i} + \beta_2 D_{2i} + u_i$$

where    $Y$ = literacy rate (percent)

$$D_1 = \text{gender} = \begin{cases} 1 & \text{if female} \\ 0 & \text{otherwise} \end{cases}$$

$$D_2 = \text{Area of residence} = \begin{cases} 1 & \text{if urban} \\ 0 & \text{otherwise} \end{cases}$$

... data of ... States of India for 2000-01 the following results we obtained (by OLS method)

$$\hat{Y}_i = \ ... \ - \ ... D_1 + \ 6.00 \ D_2$$

SE = (1.82)   (2.10)   [2.10]

$t$ = (41.65) (-7.77)  [7.62]

In this regression model there are two dummy variables. The regression ... show that the mean literacy rate is about 75.82 percent. Compared with this, the average literacy rate for female is lower by about 16 ... percent, for an actual average literacy rate of (75.82 − 18.32) = 59.50 percent.

By contrast for those who live in the urban area the literacy rates is higher by about 6 percent, for an actual average literacy rate of 75.82 + 6 = ... percent.

**Example 3.14.** This example shows regression with a mixture of quantitative and qualitative regressors. We consider the following model

Let    $Y$ = Average consumption expenditure (₹) per person per month in ... State

Let    $X$ = Average household size (the number of persons in State)

$$D_1 = \begin{cases} 1 & \text{if the State is in the Eastern region of India} \\ 0 & \text{otherwise} \end{cases}$$

$$D_2 = \begin{cases} 1 & \text{if the State is in the North-West central region of the country} \\ 0 & \text{otherwise} \end{cases}$$

The above equation is fitted with the help of the data on Household consumer expenditure in India 2000-01 and obtained

$$\hat{Y} = 2454.72 - 16.06 D_1 - ... D_2 - 44.77 X$$

SE =   503.00   145.72   160.04   28.40

$t$ =    (4.86)   (0.11)   (1.50)   (2.3)

$R^2 = 5.97$

These results suggest that other things remaining the same, as household size goes up by one person, on an average the percapita consumption expenditure goes down by about ₹ 44.72

## 114 A Brief Outline on Qualitative Response Regression Models

In all the regression models that we have considered [...] we [...] assumed that the regressand [...] whereas the explanatory [...] quantitative or a mixture thereof.

[...] we may also [...] let us consider models in which the [...] qualitative in nature. The practical importance of such models [...] in various areas of social sciences and medical research.

For example, we like to study the labour force participation [...] whether a person is either in the labour force or not. [...] hence the response variable or regressand can take [...] two values [...] a person is in the labour force and 0 if he is not. In other words the [...] is a binary or dichotomous variable.

In qualitative regression models where the regressand [...] objective is to find the probability of something happening [...] Hence qualitative response regression models are often known as probability [...]

There are four approaches to developing a probability model for a binary response variable where the regressand itself is qualitative in nature. These are

1 The linear probability model (LPM)
2 The logit model
3 The probit model
4 The tobit model

Because of its comparative simplicity, and because it can be estimated by ordinary least square (OLS) we first consider the linear probability model — LPM

### The Linear Probability Model (LPM)

We consider a two variable regression model

$$Y_i = \alpha + \beta X_i + u_i \qquad \text{(1) where}$$

$X$ = family income $Y$ = a binary variable

i.e.,
$$Y = \begin{cases} 1 & \text{if the family owns a house} \\ 0 & \text{if it does not own a house} \end{cases}$$

Model (1) looks like a typical linear regression model but because the regressand is binary it is called a **linear probability model (LPM)** This is because the conditional expectation of $Y_i$ given $X_i$, $E\left(Y_i/X_i\right)$, can be interpreted as the conditional probability that the event will occur given $X_i$, that is, $P(Y_i = 1 \mid X_i)$. Thus, in our example $E\left(Y_i/X_i\right)$ gives the probability of a family owning a house and whose income is the given amount $X_i$.

The justification of the name LPM for models like equation (1) can be seen as follows. Assuming $E(u_i) = 0$, as usual we obtain

$$E\left(Y_i/X_i\right) = \alpha + \beta X_i \qquad (2)$$

Now ... probability that Y ... that is the event ... its ... probability that ... of that is the event does not ... the ...

| Y | probability |
|---|---|
| 1 | $p$ |
| | $p$ |

This shows that $Y$ follows a **Bernoulli probability distribution**. Now by definition of mathematical expectation we obtain $E(Y) = 0 (1-p) + 1 \cdot p = p$　　(3)

Now comparing equation (3) with equation (2)

we can equate $E(Y_i) = \alpha + \beta X_i = p$　　(4)

but it is in fact the conditional probability of ... Since the probability $p$ must lie between ... and ... we have the restriction,

$$0 \leq E\left(\frac{Y_i}{X_i}\right) \leq 1 \quad \text{____ (5)}$$

From the above explanation it would seem that OLS can be easily extended to binary dependent variable regression models. So we may assume that there is nothing new here. But this is not the case because the LPM poses several problems which are as follows:

### (i) Non-Normality of the Disturbances $u_i$

Although OLS does not require the disturbances $u_i$ to be normally distributed, we assumed them to be so distributed for the purpose of statistical inference. But the assumption of normality for $u_i$ is not tenable for the LPMs because like $Y$ the disturbances $u_i$ also take only two values, that is, they also follow the Bernoulli distribution.

This can be seen clearly if we write equation (1) as $u_i = Y_i - \alpha - \beta X_i$　　(6) the probability distribution of $u_i$ is

| | $u_i$ | probability |
|---|---|---|
| when $Y = 1$ | $1 - \alpha - \beta X_i$ | $p$ |
| when $Y = 0$ | $-\alpha - \beta X_i$ | $1 - p_i$ |

　　(7)

Obvious $u_i$, $u$ cannot be assumed to be normally distributed, they follow the Bernoulli distribution. But the non-fulfilment of the normality assumption may not be so critical as it appears because we know that the OLS point estimates will remain unbiased. Besides, as the sample size increases indefinitely, statistical theory shows that OLS estimators tend to be normally distributed generally. As a result, in large samples the statistical inference of the LPM will follow the usual OLS procedure under the normality assumption.

### (ii) Heteroscedastic variances of the Disturbances

Even if $E(u_i) = 0$ and $cov(u_i, u_j) = 0$ for $i \neq j$ (i.e. no serial correlation), it can no longer be maintained that in the LPM the disturbances are homoscedastic. This is

[text heavily faded and illegible in upper portion]

... one way to restore the ... 

transform the model by dividing it through by ...

$$\frac{Y_i}{\sqrt{w_i}} = \frac{\alpha}{\sqrt{w_i}} + \beta \frac{X_i}{\sqrt{w_i}} + \frac{u_i}{\sqrt{w_i}}$$

i.e.

We can now apply OLS in this model, called **Weighted Least Square (WLS)** method with $w$ serving as weights.

In theory, what we have just described is fine. But in practice ... the ... is unknown, hence the weights $w_i$ are unknown. To estimate $w$ we can use the following two-step procedure.

**Step 1 :** We can run the OLS on regression ... despite the heteroscedasticity problem and obtain

$\hat{Y} = $ estimate of true $E\left(\frac{Y_i}{X_i}\right)$. Then obtain $w_i = \hat{Y}(1 - \hat{Y})$, the estimate of $w_i$.

**Step 2** We can use the estimated $w_i$ to transform the data shown in equation (9) and estimate the transformed equation by OLS (i.e. weighted least squares ...

**(III) Non fulfilment of $0 \leq E_i \left(\frac{Y_i}{X_i}\right) < 1$**

Since $E\left(\frac{Y}{X}\right)$ in the LPM measures the conditional probability of the event Y occurring, given X it must necessarily he between 0 an 1. Although this is true a priori there is no guarantee that $\hat{Y_i}$ the estimators of $E\left(\frac{Y_i}{X_i}\right)$ will necessarily fulfil this restriction, and this is the real problem with the OLS estimation of the LPM. This happens because OLS does not take into account the restriction that $0 < E(Y_i) \leq 1$. There

are in seaxx ... firming to whether the estimated ... is between zero and one, if the estimate lie ... he use a OLS method ... all but will ... it ... because if ... some are less than ... negative ... it is an error in the ... these also ... they are greater than ... then it ... indicate the ... or it that there is ... be because of estimating technique but we guarantee the the estimate ... the probabilities ... will lie between zero and one ... in ... and all of the estimated probabilities will ... exceed lie between the limits ... it lies 0, 1.

### (b) Questionable value of $R^2$ as a Measure of Goodness of Fit

The conventionally computed $R^2$ is of limited value in the dichotomous ... corresponding to given ... is either 0 or 1. Therefore all the ... values will either lie along the $X$ axis when $Y = 0$, or along the line ... corresponding to ... when ...

As a result the conventionally computed $R^2$ value is to be much lower than ... In most practical applications the $R^2$ ranges between ... ... ... $R^2$ ... exceeds 0.8, when the predicted $Y$ values are close to either 0 or ...

**Example 3.15.** The following table, Table 3.8, gives invented data on house ownership ($Y = 1$ owner a house, 0 = does not own a house) and income $X$ (in thousands of dollars) for 40 families.

### Table 3.8

| Family | Y | X | Family | Y | X | Family | Y | X |
|--------|---|---|--------|---|---|--------|---|---|
| 1 | 0 | 8 | 15 | 0 | 8 | 29 | 0 | |
| 2 | | 16 | 16 | 1 | 19 | 30 | | 11 |
| 3 | | 18 | 17 | | 8 | 31 | | |
| 4 | 0 | 11 | 18 | 0 | 10 | 32 | 0 | 8 |
| 5 | 0 | 12 | 19 | 0 | 8 | 33 | | 21 |
| 6 | | 19 | 20 | 1 | 18 | 34 | | 20 |
| 7 | | 30 | 21 | 1 | 22 | 35 | 0 | |
| 8 | 0 | 13 | 22 | 1 | 16 | 36 | 0 | 8 |
| 9 | 0 | 9 | 23 | 0 | 2 | 37 | 1 | 17 |
| 10 | 0 | 10 | 24 | 0 | 11 | 38 | | 16 |
| 11 | | 7 | 25 | 1 | 16 | 39 | 0 | |
| 12 | 1 | 18 | 26 | 0 | | 40 | | 7 |
| 13 | 0 | 14 | 27 | 1 | 37 | | | |
| 14 | | 30 | 28 | | 8 | | | |

From these data the estimated LPM (by OLS method) is given below

$$\hat{Y}_i = -0.9457 + 0.1021 X_i$$

SE    (0.1228)    (0.0082)

t = ...    (2.4913),    $R^2 = 0.8048$

From the estimated LPM the intercept of -0.9457 gives the 'probability' that a family with zero income will own a house. Since this value is negative and since probability

$$\hat{Y}_i \; X=12) = \; 0.9457 + 12 = 0.1032 = 0.2794$$

... the probability that ...

We ...

We can make the probabilities ...

estimated probabilities will be negative and some will exceed 1

... for instance when $X = 8$, ( $\hat{Y}_i$ , $X = 8$) = -0.9457 + 8 = 0.102

Similarly, when $X = 20$, ( $\hat{Y}_i$ , $X = 20$) = 0.9457 - 20 + 0.102

this means that although ...

... not be necessarily positive or less than ... This ...

... coefficient more when the dependent variable ...

... even if the estimated Y are all positive and ... the ...

problem of heteroscedasticity. () it is ... we have to ... app ...

make.

# EXERCISE

1. In a multiple linear regression model, there are ... 
   (ii) ...    ii. why do we insert the random disturbance term

2. State the assumptions about ... of a classical linear regression model ... if the model assumes the form $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$, $i = 1, 2$

3. In a three variable linear regression model of the form ... 
   ... ii. how can you estimate the regression parameters ... method?

4. Describe briefly the method of least squares used in estimating the regression parameters relating to a three variable linear regression model.

5. State and prove the properties of the least squares estimators relating to a three variable linear regression model (CLRM)

6. Show that in a three variable classical linear regression model 
   $u_{ii} = 2$    a. the estimated parameters coefficients are unbiased

7. In a CLRM, $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$, ... a. express the estimated regression coefficients $\beta_1$ and $\beta_2$ in terms of variables and coefficient of deviations.

8. Describe the variances and covariances of the regression parameters in the model 
   $Y = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$, $i = 1, 2, \ldots, n$.

9. State and prove GAUSS-MARKOV THEOREM in terms of a three variable linear regression model of the form $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$, $i = 1, 2, \ldots n$

10. What is meant by ... regression line ...

11. How can you determine ... draw the ... regression line ... $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + e_i, i = 1, 2, \ldots, n$?

12. In terms of a three variable linear regression model, show that $\sum \ldots$ ... unbiased estimate ... random error term ...

13. When is a multiple least squared ... multiple linear regression model ... find?

14. Show that in a three variable linear regression model $M_0$ ... is $\sum \ldots$ unbiased estimate of ...

15. Show that the least square estimator of the model ... that
    (i) $\hat\beta_0 = \ldots$
    (ii) $\ldots$
    where $E[\hat\beta_i] = \beta_i$, ... $= \beta_i + \beta_i$.
    $var(\hat\beta_0) = \ldots$
    $var(\hat\beta_i) = \ldots$

16. Show that MLE of ... in a three variable linear regression model is not an unbiased estimator of $\sigma_u$ but a consistent ... biased estimate ...

17. Describe the testing procedure of the regression ... regression ... parameters ... model $Y = \beta_0 + \beta_1 T ... \beta_2 X ...$ ... $= 1, \ldots n$.

18. What is meant by goodness of fit of a ... three linear regression model.

19. What is meant by multiple coefficient of determination. Derive the formula of multiple coefficient of determination of a series in a three variable linear regression model.

20. In terms of ... ... ...

21. In terms of a linear regression ...

22. In terms of a three-variable linear regression ...

the coefficient $R$ ...

where ...

23. Define multiple coefficient of determination $R^2$ ...
relation between $R$ and $R$ when there is ...

24. What is adjusted $R^2$? Why is it used ...

25. Define partial correlation coefficient ...
in a three-variable linear regression model.

26. Establish the relation among multiple coefficient of determination $R^2$, ...
correlation coefficients and partial correlation ... three-variable linear regression model.

27. How can you formally write the regression model ...
$= \mu_i X_{3i} + u_i$, where $u_i$ is ... $+ 2 ...$ +) mention all the properties of CLRM.

28. How can you use the analysis of variance ... linear regression ...

29. What is the meaning of the beta production ... How are ... expressed in a three-variable linear regression model. Distinguish between partial and ... ... probation in this regard.

30. Establish the relation between $R^2$ and $F$ in terms of a three-variable linear regression model. What would be the value of $F$ when $R^2 = 0$?

31. What is an ANOVA Table? How can you construct an ANOVA table in terms of $R^2$?

32. The Cobb–Douglas production function in its stochastic form is given by ...
$Y_i = \beta_1 X_2^{\beta_2} e^{u_i}$, where $Y$ = output, $X_2$ = labour input ... ... stochastic disturbance term, $\beta_0$ = constant technological parameter. How ... estimate the regression parameters by applying OLS method. What would ... ... ... for this function?

33. What do you mean by indicator variables, ... variables, qualitative variables, dummy variables, binary variables? Give an example.

34. What do you mean by dummy variables? How can you incorporate these variables in the regression model. What is meant by dummy ... trap?

35. What do you mean by dummy variables? Explain some ... the uses of dummy variables in applied economic research.

**36.** What are the dummy variables? Construct a model where dummy variables are used as proxies for the dependent variable.

**37.** What are the dummy variables? Construct a model where dummy variables are used as proxies of numerical factors.

**38.** What are the dummy variables? Construct a model where dummy variables are measuring the shift of a function in time.

**39.** What do you mean by qualitative response regression models? How are the qualitative regression models made useful to applied econometric research?

**40.** The following were obtained from 15 sets of observations on $Y$ and $X$ ... Estimate the regression of $Y$ on $X$ and comment on the results. ... that the coefficient of $X$ is zero.

**41.** Consider the following regression model in deviation form $Y_i = \beta_1 + \beta_2 X_{2i} + u_i$

sample data ...

Then:

    Compute OLS estimates of $\beta_1$, $\beta_2$ and $R^2$.
    Test the hypothesis $H_0: \beta_2 = 0$ against $H_1: \beta_2 \neq 0$.
    Test the hypothesis $H_0: \beta_1 = \beta_2 = 0$ against $H_1: \beta_1 \neq 0, \beta_2 \neq 0$.
    Test the hypothesis $H_0: \beta_1 = \beta_2$ against $H_1: \beta_1 \neq \beta_2$.

**42.** A production function model is specified as $Y = \beta_0 + \beta_1 L + \beta_2 K$ where $Y$ = log output, $L$ = log labour input and $K$ = log capital input. The data refer to a sample of 15 items and observations are measured as deviations from the sample means.

$\Sigma ... \quad \Sigma ... = X, \Sigma ... , \Sigma ... , 10, \Sigma ... = N, \Sigma ... , 13$

    Estimate $\beta_1$, $\beta_2$ and their standard errors
    (ii) Find $R^2$ and adjusted $R^2$
    (iii) Test the hypothesis that $\beta_1 = \beta_2$
    (iv) Suppose now that you wish to impose the apriori restriction that $\beta_1 + \beta_2 = 1$. What is the least squares estimate of $\beta_1$ and its standard error? What is the value of $R^2$ in this case? Compare these results with those obtained in ... and estimates.

**43.** The following table shows 10 sets of values of three variables $Y$ (dependent variable $Y$ and $X_2$ (two independent variables):

| $Y$ | 3.5 | 4 | 5 | 6 | 9 | 8 | 7 | 2 | 4 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $X_1$ | 5 | 20 | 38 | 42 | 50 | 24 | 65 | 71 | 85 | 90 |
| $X_2$ | 16 | 13 | 10 | 7 | 7 | 3 | 4 | 3 | 35 | 7 |

(i) Consider a model of the form $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u$. Find the least squares regression equation of $Y$ on $X_1$ and $X_2$.

(ii) Compute the coefficient of multiple determination and the standard errors of the estimated parameters and conduct tests of significance.

(iii) Construct 95 percent confidence intervals for the population parameters and ...

(iv) Find the explained and unexplained variation in $Y$.

... coefficients

... the following ... estimates and test ... significance

... A. D ...

b. that ...

... the multiple ...

... obtain 95% confidence ...

The following table shows the value of imports ... measured in arbitrary units and the per ... twelve-year period for a certain country:

| Year | 1989 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | | | |
|------|------|------|------|------|------|------|------|------|--|--|--|
| ... | 4. | | 37 | 44 | | 16 | | | | | |
| 720 | | | 243 | 405 | 253 | 354 | 321 | 376 | 375 | 382 | 344 |
| | 42 | 118 | 140 | 228 | 349 | 145 | 150 | 148 | 153 | 155 | 153 |

i. Estimate the import function $Y = \beta_0 + \beta_1 X_1 + u$

ii. What is the economic meaning of your estimates.

iii. Construct tests of significance for the regression estimate ... and ... significance

iv. Compute $R^2$ and adjusted $R^2$

iii. The following table includes the output ... for ... firms of the chemical industry:

| Firm | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|------|---|---|---|---|---|---|---|
| Y (000 tons) | 60 | 120 | 190 | 240 | 300 | 360 | 390 | 430 |
| L (hours) | 300 | 1200 | 430 | 500 | | | | |
| K (machine hours) | 300 | 400 | 470 | 400 | 510 | 440 | 600 | 670 |
| | 9 | 10 | 11 | 12 | 13 | 14 | 15 | |
| | 440 | 490 | 500 | 520 | 540 | 410 | 350 | |
| | 1800 | 1750 | 1950 | 1960 | 1610 | 1900 | 150 | |
| | 620 | 640 | 850 | 900 | 800 | 900 | 800 | |

i. Fit a Cobb-Douglas production function to the above data $Y = b_0 L^a K^b u$

ii. Construct appropriate tests of significance of the parameter estimates at ... and ... % levels of significance

iii. What are the marginal and average productivities of the factors $L$ and $K$?

iv. What do your results suggest regarding the returns to scale.

**47** The following table shows the price index of durables, the average yearly household expenditure on durables of a typical household of a country

| Year | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 |
|------|------|------|------|------|------|------|------|------|------|
| Expenditure on durables $Y$ (in £) | | 5 | 10 | 115 | | 40 | 80 | 75 | 45 |
| Household income ($X$) (£) | 835 | 2000 | 030 | 2040 | | | 905 | 915 | |
| Price index ($Y$) | 100 | 05 | 95 | 05 | 95 | | | 15 | 5 |

(i) Fit a regression line to the function $Y = \beta_0 + \beta_1 X + \beta_2 Z + u$

(ii) Test your results by using the Analysis of variance table

**48.** The following table shows the consumption of tobacco manufactures, consumer's income and the price of tobacco manufactures for France during 1950s.

| Year | Consumption (million tons) $D$ | Income (million francs) $I$ | Price of tobacco (francs per kg) $P$ |
|------|------|------|------|
| 1950 | 59,190 | 6,200 | 2 46 |
| 1951 | 65,450 | 91,700 | 24 44 |
| 1952 | 62,360 | 98,700 | 19 7 |
| 1953 | 64,700 | 111,600 | 32 46 |
| 1954 | 67,400 | 119,800 | 31 49 |
| 1955 | 64,440 | 129,200 | 34 44 |
| 1956 | 68,000 | 143,400 | 35 30 |
| 1957 | 72,400 | 159,600 | 38 70 |
| 1958 | 75,710 | 180,000 | 39 43 |
| 1959 | 78,640 | 191,000 | 44 68 |

(i) Fit a linear regression $D = \beta_0 + \beta_1 P + \beta_2 I + u$ and a non-linear function of the constant elasticity type $D = \beta_0 P^{\beta_1} Y^{\beta_2} u$.

(ii) Conduct tests of significance using the analysis of variance table

(iii) Compute the price and income elasticities of the two functions

**49.** In a multiple regression equation $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + u$ explain how you would test the joint hypothesis $\beta_1 = \beta_2$ and $\beta_2 = 1$

**50.** Consider the following regression model $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$ where $u_i \sim N(0, \sigma_u^2)$. The following data set are given below

| Y | 4 | 7 | 3 | 9 | 7 |
|---|---|---|---|---|---|
| $X_1$ | 2 | 3 | 1 | 5 | 9 |
| $X_2$ | 5 | 3 | 2 | 1 | |

(i) Estimate $\beta_0$, $\beta_1$ and $\beta_2$

(ii) Find out $\mathrm{var}(\hat{\beta}_0)$, $\mathrm{var}(\hat{\beta}_1)$, $\mathrm{var}(\hat{\beta}_2)$ and $\mathrm{cov}(\hat{\beta}_1 \hat{\beta}_2)$?

(i) Find R and ...

(ii) Write the regression ...

(iii) ...

(iv) ... test $H_0$ ... against ...

(v) ...

(vi) ...

(ix) Construct ANOVA and ...

(x) Find out the point prediction ...

(xi) Construct 95% confidence interval ...

(xii) Test whether the effect ...

55. Eight students made the following ... certain subject. For the linear regression ... score $X_1$ and test score $X_2$.

| Students | | | | | | | |
|---|---|---|---|---|---|---|---|
| Pretest score $X_1$ | 43 | | | 28 | | | |
| Test score $X_2$ | | 29 | 56 | | 26 | | |
| Final score $Y$ | | 34 | 35 | | 14 | | |

Assume a linear regression equation of the form

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

(i) Find $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$.

(ii) Find $\text{var}(\hat{\beta}_0)$, $\text{var}(\hat{\beta}_1)$, $\text{var}(\hat{\beta}_2)$ and $\text{cov}(\hat{\beta}_1, \hat{\beta}_2)$.

(iii) Find $R^2$ and adjusted $R^2$.

(iv) Write the regression results in the summary form.

(v) Test whether $\beta_0$, $\beta_1$ and $\beta_2$ are significant or not at 5% level of significance.

(vi) Find out the point prediction of $Y$ when $X_1 = 44$ and $X_2 = 62$.

(vii) Construct 95% confidence interval of $\pi_1$.

(viii) Construct 95% confidence interval of the prediction.

51. Following exercise-50, test the regression results by using Analysis of variance table.

52. Following exercise 51, test the regression results by using Analysis of variance table.

54. The following table shows monthly income $Y$ (in ₹), monthly savings $X_1$ (in ₹), and age $X_2$ of 10 persons.

| $Y$ | 2000 | 8000 | 7000 | 2000 | 10000 | 25000 | 35000 | 40000 | 30000 | 40000 |
|---|---|---|---|---|---|---|---|---|---|---|
| $X_1$ | 0.000 | 20.000 | 15.000 | 30.000 | 30.000 | 45.000 | 50.000 | 55.000 | 60.000 | 70.000 |
| $X_2$ | 22 | 24 | 40 | 35 | 41 | 43 | 46 | 50 | 57 | 60 |

Assume a linear regression equation of the form $Y = \beta_0 + \beta_1 ...$

and assume $X_2 = 1$ ... $= 0$ if $< 40$ years of age
$= 1$ if $\geq 40$ years of age

(i) Find $\beta_0$, $\beta_1$ and $\beta_2$

(ii) Find $R^2$ and adjusted $R^2$

(iii) Write the regression results in the standard form

(iv) Test whether $\beta_0$, $\beta_1$ and $\beta_2$ are significant or not at 5% level of significance

(v) Obtain 95% confidence interval of $\beta_2$

**55.** The following regression was estimated from ... (the figures in parentheses)

$$\hat{Y} = 0.06 + 1.41 X_1 + 6.64 X_2 + ... \qquad R^2 = ...$$
$$(3.7) \quad (0.27) \quad (1.37) \quad (5.40) \quad (3.17)$$

where $X_{i}$ ... the $i$th quarter and 0 otherwise. Explain the implied pattern of seasonal variation and interpret the result.

**56.** You are given the following regression results

$$\hat{Y}_i = 16.890 - 2972.5 X_{2i} \qquad R^2 = 0.6149$$
$$\phantom{\hat{Y}_i =}(8.5 \cdot 52) \qquad (476.9)$$

$$\hat{Y}_i = 897.4 - 2.14 X_i \qquad t = 25.3, r \qquad R^2 = 0.7906$$
$$(3.7054) \qquad (0.6070) \qquad (2.9712)$$

Can you find out the sample size underlying these results?

*Hint:* Use the relationship among $R^2$, $F$ and $t$ values.

**57.** From the data for 46 States in the United States for a given year the following regression results were obtained

$$\log \hat{Y} = 4.30 - 1.34 \log P + 0.17 \log Y$$
$$SE \qquad (0.91) \quad (0.32) \qquad (0.20) \qquad\qquad R^2 = 0.27$$

where $Y$ = state of consumption of a commodity per year
$P$ = real price per unit of the commodity
$Y$ = per capita real disposable income

(i) What is the elasticity of demand for the commodity with respect to price? Is it statistically significant? If so, is it statistically different from 1?

(ii) What is the income elasticity of demand for the commodity? Is it statistically significant?

(iii) How would you retrieve $R^2$ from $\bar{R}^2$ given above?

**58.** From a sample of 209 firms the following regression results were obtained

$$\log(\text{salary}) = 4.32 + 0.280 \log(\text{sales}) - 0.0174 \text{ roe} + 0.00024 \text{ ros} \qquad R^2 = 0.29$$
$$SE \qquad (0.32) \quad (0.035) \qquad\qquad (0.0041) \qquad (0.00054)$$

where salary = salary of CEO

sales = annual firm sales

roe = return on equity in percent

ros = return on firm's sales

Interpret the regression ... and ... ... ... ... ...
... have about the signs of the various coefficients

ii. Which ... ... ... ... are ... ... ... ... ... ...

iii. What is the ... statistical significance of ... regression

iv. Can you interpret the coefficient of ... ... ... ... ... ... ... or why not?

**59.** Consider the following wage determination ... ... ... ... ... ... ... period 1951-1969

$$W = 8.582 + 0.364(PF)_t + 0.004(PF)_{t-1} - 2.560U \qquad R^2 = ...$$
$$t = (1.29) \qquad (4.180) \qquad (0.72) \qquad ...$$

where W = wages and salaries per employee

PF = price of final output at factor cost

U = unemployment in Great Britain as a percentage ... ... ... employees in Great Britain

t = time

i. Interpret the regression equation

ii. Are the estimated coefficients individually significant?

iii. What is the rationale for the introduction of PF?

iv. How would you compute elasticity of wages and salaries per employee ... ... ... to unemployment rate U?

**60.** Consider the following data set

| Y | 3 | 4 | 5 |
|---|---|---|---|
| X₁ | 1 | 2 | 3 |
| X₂ | 2 | 4 | 7 |

(in LaTeX:)

| $Y$ | 3 | 4 | 5 |
| $X_1$ | 1 | 2 | 3 |
| $X_2$ | 2 | 4 | 7 |

Based on these data, estimate the following regressions

$$Y_i = \alpha_0 + \alpha_1 X_{1i} + u_{1i} \qquad (1)$$

$$Y_i = \lambda_0 + \lambda_2 X_{2i} + u_{2i} \qquad (2)$$

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i \qquad (3)$$

Estimate the regression coefficients in each case

ii. Is $\alpha_1 = \beta_1$? Why or why not?

(iii) Is $\lambda_2 = \beta_2$? Why or why not?

What important conclusions do you draw from this exercise?

**61.** From the following data estimate the partial regression coefficients, their standard errors, and the adjusted and unadjusted $R^2$ values

$$\bar{Y} = 367.693, \quad \bar{X}_1 = 402.760, \quad \bar{X}_2 = 8.0$$

$$\Sigma(Y_i - \bar{Y})^2 = 66042.269 \quad \Sigma_i \lambda_{1i} \quad X_1 = 84855.096$$

$\Sigma(x_1 - \bar{x})$ ... $\Sigma(x_1 - \bar{x})(x_2 - \bar{x}) = 14772.146$

$\Sigma(x_1 \bar{x})$, $\bar{x} = 4250.900$ $\Sigma(x_{2} - \bar{x})(x_{2}) = 4196.00$ $n = 5$

**62.** Is it possible to obtain the following from a set of data ?

(i) $r_{23} = 0.9$, $r_{12} = 0.2$, $r_{13} = 0.4$

(ii) $r_{12} = 0.6$, $r_{23} = -0.9$, $r_{31} = 0.5$

(iii) $r_{31} = 0.01$, $r_{12} = 0.66$, $r_{23} = 0.7$

**63.** The regressed child mortality (C.M) on percapita GNP (PGNP) and the female literacy rate (FLR) for a sample of 64 countries is given below

$$\widehat{CM}_i = 263.64 - 0.0056\ PGNP_i - 2.23 + FLR_i$$

$$SE \quad (11.5932) \qquad (0.0019) \qquad (0.2099)$$

$$R^2 = 0.7077, \quad \bar{R}^2 = 0.6981$$

(i) Interpret the regression results

(ii) What about the statistical significance of the observed results ?

(iii) Is the coefficient of PGNP of −0.0056 statistically significant ?

(iv) Is the coefficient of FLR of −2.2316 statistically significant ?

(v) Are both the coefficients statistically significant jointly ?

# 4

# Violations of Classical Assumptions: The Problems of Heteroscedasticity, Autocorrelation and Multicollinearity

## 4.1 Introduction

Let us consider a two-variable linear regression model

$$Y_i = \alpha + \beta X_i + u_i$$

If the model is a classical linear regression model, then it satisfies the following assumptions:

(i) $u$ is a random variable

(ii) $E(u_i) = 0$ for each $i$

(iii) $E(u_i^2) = \sigma_u^2$ or $\sigma^2$ (constant)

(iv) $Cov(u_i, u_j) = E(u_i, u_j) = 0$ for $i \neq j$

(v) $X$ is non-stochastic or non-random

We now put special consideration of assumptions (iii) and (iv). Assumption $E(u_i^2) = \sigma^2$ means that the variance of the disturbance terms is constant which is the probability distribution of the disturbance term does not vary. This is the feature of homogeneity of variance (or constant variance). It shows, however, that it may be the case, however, that all of the disturbance terms do not have the same variance. This condition of non-constant variance or non-homogeneous variance is known as *heteroscedasticity*. Thus we say that $u$ are heteroscedastic when

$$var(u_i) = \sigma_u^2 \text{ (or } \sigma^2\text{), a constant value] but var } u_i = \sigma_i^2 \text{, a variable variance]}$$

Also we should not assume that each disturbance term has the same expected value equal to zero i.e. $E(u_i) = 0$

If all the disturbance terms have expected value zero and same variance $\sigma^2$ then we say that all the disturbance terms are identically distributed.

If (iv) is also satisfied, it means that the different disturbance terms are independent of each other. So, when all (ii), (iii) and (iv) are satisfied, then we can say that the different disturbance terms are identically and independently distributed.

If the disturbance term varies from observation to observation, the different disturbance terms are not identically distributed.

Here $var(u_i) = E(u_i^2) = \sigma_u^2 \neq \sigma_u^2$ when $i \neq j$

This is the problem of *Heteroscedasticity*. The CLRM assumes that variance of the disturbance term is constant and if it is not constant, then the problem of

Autocorrelation is a special case of correlation. It refers to the relationship between the successive values of the same variable in different periods while correlation refers to the relationship between two or more different variables.

Thus, in the presence of autocorrelation the different disturbance terms are not independent of each other.

So, both the assumptions of CLRM are violated there is the problem of heteroscedasticity as well as the problem of autocorrelation.

## 4.2 Matrix Representation of Autocorrelation and Heteroscedasticity

Let us consider the dispersion matrix or variance-covariance matrix of the disturbance vector

### 4.3 Consequences of the Problem of Autocorrelation and Heteroscedasticity

### 4.4 Consequences of the Problem of Heteroscedasticity

The OLS estimator of $\beta$ is denoted by $\hat{\beta}$

### 4.3 Method for Estimating Regression Parameters in the Presence of the Problem of Heteroscedasticity

$$var(\hat{\beta}) = \frac{\sigma^2 \sum_{i=1}^{n} x_i^2}{\left(\sum_{i=1}^{n} x_i^2\right)^2} \left[ \sum_{i=1}^{n} x_i^2 = \sigma_i^2 x_i \right]$$

with $ \ldots \ldots $

The WLS estimator ... the parameter $\beta$ denotes to $R^2$ ... a lower variance than the ... variance of $\beta$

In this sense WLS method is more appropriate to a model which is subject to the problem of heteroscedasticity.

(Note: Heteroscedasticity may also be of the form $\sigma_i^2 \ldots \ldots \ldots$)

## 4.6. Tests for Heteroscedasticity

There are three important tests for heteroscedasticity.
(i) Spearman's Rank correlation test
(ii) Goldfeld and Quandt test.
(iii) Glejser's test.

All these test criteria are based on the OLS residual term.

Let us consider the following model: $Y = \alpha + \beta X + u$ where ...

... of the given observations on $Y$ and ... and apply on the ... ... over $y$ and $\hat{y}$ and then on the get ... $ = \hat{y}^2 \ldots$ ...

... to obtain the ... $\hat{u}$ Random error but after these steps ...

### Spearman's Rank Correlation Test

... for large or small samples ...
... to compute the rank correlation coefficient between ...

... rank coefficient is computed by the formula

$$ r = 1 - \frac{6 \sum d_i^2}{\ldots} $$

... Difference between the ranks of corresponding pairs of $X$ and ... and ... number in the sample.

... to test the null hypothesis that the rank correlation coefficient ...
the alternative hypothesis that it is not equal to zero ... the null ... and the alternative hypothesis. $H_0 : \ldots = 0$

... $t$-ratio test statistic is then given by ...

$$ t = \ldots \ldots $$

... has $t$-distribution with $n - 2$ degrees of freedom.
... If a large sample the null hypothesis $H_0 : \ldots$ it will be ... rejected ... ... with ... $r$ ... If ... is ... good to the actual ...
... If the null hypothesis is accepted then there is no problem of ... . But if it is rejected then there will be the problem of ... ...

### Goldfeld and Quandt Test

... applicable to large samples. For this we have to ... consider the ...
... order the observations according to the magnitude of the ... variable $X$.
... a series, an arbitrary number ... and then a certain number of central ... ... from the middle. The remaining observations ...
... are then divided into two equal parts, each part therefore has ... observations. One part includes the small values of $X$ while the other part ... the large values of $X$.
... structure regression lines by OLS procedure to each part and obtain the ... residuals from each of them.
... assume the sum of squared residuals from the sample of low values ... and ... denote the same from the sample of large values of ... Then we calculate

### 4.4 Mean, Variance and Covariance of the Autocorrelated Disturbance Variable

**1. Mean of the autocorrelated $u$'s**

**2. Variance of the autocorrelated $u$'s**

**3. Covariance of the autocorrelated $u$'s**

### Consequences of Autocorrelation

$$n + \sum \dots \sum \dots$$

for large values of $n$ the large values $\sum \dots$ and $\sum \dots$ are approximately equal

$$Then \quad \dots$$

But $\dots$ where $\rho$ is the first estimate of $\dots$ in the model $\dots$

This shows that $d$ lies between 0 and 4.

**Firstly,** if there is no autocorrelation $\rho = 0$ and $d = 2$. Then if from the sample the final $d$ (observed) $= d = 2$, we accept that there is no autocorrelation in the function.

**Secondly,** if $\rho = +1$, $d = 0$ and we have perfect positive autocorrelation.

**Thirdly,** if $\rho = -1$, $d = 4$ we have perfect negative autocorrelation. Therefore if $2 < d < 4$ there is some degree of negative autocorrelation, which is stronger if the value of $d$ is higher.

If there is no autocorrelation $\rho$ should be zero and $d = 2$.

Thus we have to set the null hypothesis $H_0$: $d = 2$ against the alternative of $d \neq 2$.

Here the problem is that the exact values or distribution of the values of $d$ is not known. What Durbin and Watson have done is to specify an upper and a lower limit of $d$.

Let $d_U$ stand for the upper limit of $d$ and $d_L$ stand for the lower limit of $d$. This is shown in the diagram on the next page (Fig 4.1)



**Fig 4.1 The critical regions of d are shown**

With the help of $d_U$ and $d_L$ we have to determine whether autocorrelation exists once the values of $d$ and $d_L$ are regulated at the $5\%$ and $1\%$ level of significance.

$$d \text{ calculated} = \dots = \frac{\sum \dots}{\sum \dots} \quad \text{from the sample}$$

From this we interpret the following cases:

1. If $d_U < d < 4 - d_U$ we reject the null hypothesis of no autocorrelation and accept that there is positive autocorrelation of first order.

2. If $4 - d_L < d < 4$ we reject the null hypothesis of no autocorrelation and accept that there is negative autocorrelation of the first order.

3. If $d_U < d < 4 - d_U$ we accept the null hypothesis of no autocorrelation.

4. If $d_L < d < d_U$ or $4 - d_U < d < 4 - d_L$ the test is inconclusive.

**Limitations of Durbin-Watson Test**

There are several limitations of Durbin-Watson test.

1. There exist inconclusive regions. If the value of $d$ lies between either $d_L$ and $d_U$ or between $4 - d_U$ and $4 - d_L$, then we cannot conclude whether autocorrelation exists or not.

2. The test method is appropriate only when the nature of the autocorrelation is of first order autoregressive type. But it is not appropriate when autocorrelation is of higher order and non-linear type.

3. If there is any lagged variable as independent variable in the model, then Durbin-Watson statistic $d$ is inappropriate in testing for autocorrelation.

### 4.10.1 Von Neumann Ratio Method of Testing Autocorrelation

This ... the ratio of the variance of the first differences of ...

$$\delta^2 = \frac{\sum_{i=2}^{n}(z_i - z_{i-1})^2}{n-1}$$

The Von Neumann ratio is applicable for ...

$$\frac{\delta^2}{S^2} = \frac{\sum_{i=2}^{n}(z_i - z_{i-1})^2 / (n-1)}{\sum(z_i - \bar{z})^2 / n}$$

This test statistic is used for measuring the existence of autocorrelation.

This procedure is, however, not applicable for testing the autocorrelation of the ...

## 4.11 Methods for Estimating Regression Parameters in the Presence of the Problem of Autocorrelation

Once the autocorrelation is detected, the appropriate remedial procedure ... OLS will estimate ... the regression data ...

Let our model be given by $Y_t = \alpha + \beta X_t + u_t$,

where $u_t = \rho u_{t-1} + \varepsilon_t$ with ...

If we take a lagged form of the model ... multiply both sides by $\rho$, we obtain

$$\rho Y_{t-1} = \rho \alpha + \rho \beta X_{t-1} + \rho u_{t-1}$$

Now subtracting ... from ... you get

$$Y_t - \rho Y_{t-1} = \alpha(1-\rho) + \beta(X_t - \rho X_{t-1}) + \varepsilon_t$$

or, $Y_t^* = \alpha^* + \beta(X_t - \rho X_{t-1}) + \varepsilon_t$

or $Y_t^* = \alpha^* + \beta X_t^* + \varepsilon_t$

where $Y_t^* = Y_t - \rho Y_{t-1}$, $X_t^* = X_t - \rho X_{t-1}$ and $\alpha^* = \alpha(1-\rho)$

Here $\varepsilon_t$ satisfies all the properties of CLRM.

It should be noted that in transforming ... the observations shall be lost because of lagging and subtracting ... from ... we can apply OLS to the transformed relation ... to obtain $\alpha^*$ and $\beta$ for our parameters $\alpha$ and $\beta$ ...

## 4.... AUTOCORRELATION AND MULTICOLLINEARITY 207

$$... + ... = \frac{1}{n}\sum ... \quad \text{and } v = \frac{u^2}{n} ...$$

... $\hat{z}$ is perfectly and linearly related to ...

... all variable appropriate of ... the variances $\alpha$, $\beta$ and $\rho$ are standard OLS formulae.

$$\beta = \frac{\sum ...}{...}$$

$$\alpha = \frac{\sum ...}{...}$$

$$\rho = \frac{\sum ...}{...} \quad \text{where } ... $$

**Method I: Using consensus estimation on $\rho$**

Sometimes on several grounds we make a ... reasonably good about the value of the ... such knowledge of ... about the ... equation ... usually we assume ... from the transformed model becomes

$$... = ... + ...$$

... method we get estimate only $\hat{\alpha}$ but $\hat{\beta}$ cannot be estimated.

If $\rho = 1$ then $\alpha$ and $\beta$ can be estimated at a time.

**Method II: Estimation of $\rho$ from the statistics.**

From the Durbin-Watson test statistic we know that $d = 2(1-\rho)$

where $d = \frac{\sum e_t e_{t-1}}{\sum e_t^2}$

Suppose that we calculate certain value of $d$ statistic $= d^*$ from sample data.

that $d^* = ...$ so that $\hat{\rho} = ... d^*$

If ... is estimated then we can estimate $\alpha$ and $\beta$ from the model.

It should be noted that $\hat{\rho}$ will not be accurate if the sample size is small. This is ... only for large samples.

The page is too faded and low-resolution to produce a reliable transcription.

## 6.9 Multicollinearity—Meaning and Sources

## 6.10 Consequences of Multicollinearity

### 6.10.1 Exact Multicollinearity and its Consequences

## 4.8. Some Illustrative Examples

We can discuss some examples where the intercorrelations between the explanatory variables are high and show the consequences on under the model

Example 8.5 ...

## 4.12 Solutions to the Problem of Multicollinearity

There are the commonly used methods to solve the problem of multicollinearity:
1. Dropping of variables
2. Using extraneous estimates
3. Ridge Regression
4. Using ratios or first differences
5. Using Principal components
6. Getting more data

## Using First Differences

The CLRM properties are

$$E(u_t) = 0, \quad E(u_t^2) = \sigma_u^2, \quad E(u_t u_{t-1}) = 0$$

$$E(\Delta u_t^2) = E[u_t - u_{t-1}]^2$$

$$= E[u_t^2 - 2u_t u_{t-1} + u_{t-1}^2]$$

$$E[u_t] = {}^2E[u_t u_{t-1}] \quad E = 0$$

$$\sigma_u^2, \; 2\sigma_u^2, \; 0 \quad \sigma_u^2 = 2\sigma_u^2$$

## EXERCISE

1. What is meant by autocorrelation? How does it arise? ...

2. What is the meaning of the term heteroscedasticity? ...

3. ...

4. Explain the possible consequences of the problem of heteroscedasticity.

5. Show that the OLS estimators will be unbiased even if there is a problem of autocorrelation in the CLRM.

6. Show that in presence of heteroscedasticity the OLS estimators will be unbiased but BLUE property may not be retained.

7. How can you estimate the regression parameters in presence of heteroscedasticity?

8. Explain the WLS method in estimating a regression model in the presence of ... of heteroscedasticity.

9. Explain briefly the different testing procedures used for testing ... of heteroscedasticity.

10. Given $Y_i = \alpha + \beta X_i + u_i$, with $E(u_i^2) = \sigma^2$, ... prove that OLS estimator ... possesses greater variance than the OLS estimates of the transformed version ...

11. Consider the model $Y_i = \alpha + \beta X_i + u_i$, ...

# 5

## Specification Analysis

### 1 Introduction

### 5.2 Diagnostic Tools Based on Least Squares Residuals

### Model Selection Criteria

### 1.1 Types of Specification Errors

...

$$\dots$$

where ...

...errors of measurement bias.

...

...education

...relative to students with ...education

...

In developing an empirical model, one is likely to commit one or more of the following specification errors:

1. Omission of a relevant variable(s)
2. Inclusion of an unnecessary (irrelevant) variable(s)
3. Adoption of the wrong functional form
4. Errors of measurement
5. Incorrect specification of the stochastic error term
6. Assumption that the error term is normally distributed

...specification errors, in the sense that we have in mind a "true" model but we do not in reality use the correct model. In the latter case we run into what may be called **model mis-specification errors**, for we do not know what the true model is to begin with.

In the present book we will concentrate mainly on the ...specification errors.

### Consequences of Model Specification Errors

...

### Underfitting a Model (Omitting a Relevant Variable)

...

$$\beta_1 + \beta_2 \dots$$

where ...

...

### 5.6.2 Tests for Omitted Variables and Incorrect Functional Form

[Text heavily faded and largely illegible]

#### (i) Examination of Residuals

[Text heavily faded and largely illegible]

$$\eta_t = \beta_0 + \beta_1 x_t + u_t$$

but researcher fits the following quadratic cost function

$$\eta_t = \alpha_0 + \alpha_1 Y + \ldots \tag{2}$$

and another researcher fits the following linear cost function

$$\eta_t = \lambda Y + u_t$$

where $Y$ = total cost, $X$ = output

[Text faded] ... the true cost function then both the researchers ... specification ...

Let us consider an example.

**Example 5.3.** We have the following set of data on output $(X)$ and total cost $(Y)$:

| Output | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Total | | | | 240 | | | | | |

[Column heavily faded and largely illegible]

$$\hat{Y} = \ldots$$

$$\hat{Y}_t = \ldots \quad R^2 = 0.9219 \quad \text{where } \ldots$$

[Text faded]

[Text faded] ... is saved in ... all the data and functions are stored in the following table ...

#### Table 5.3.

**Estimated residuals from the Linear, Quadratic and Cubic total cost functions**

| Y residuals | Linear Model | Quadratic Model | Cubic Model |
|---|---|---|---|
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

where $R^2_{...} = 0.9981$, $R^2_{...} = 0.5109$

... number of new ... number ... in the new model = 4

... we obtained $F$ value ... and ... at ... ... indicating that the model ...

... 5): (9902.1) ... (1.05)

is as expected ...

... once we have reached ... ... by the ... ... examination of the residuals ... Durbin-Watson ... value ...

... one advantage ... ... ... will ... be us specify who the ... ... possible is that the ... ... from average because knowing but a better ... ... does not help us necessarily in choosing a better alternative.

### (iv) Lagrange Multiplier (LM) Test for Adding Variables

This test is an alternative to ... (F/N) ... test for detecting specification errors in a model.

In order to explain this test procedure here also we are using the examples of cost functions. Let $Y = \beta_0 + \beta_1 X$ ... be a linear cost function and

$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3$ ... be a cubic cost function where

$Y$ = Total cost and $X$ = output

If we now compare the linear ... function with the cubic cost function, then the linear cost function will be a restricted version of cubic cost function. The restricted regression ... will ... assuming all the coefficients of the squared and cubed output terms i.e. $\beta_2$ and $\beta_3$ are equal to zero.

To test that the add test runs as follows:

i. We estimate the restricted regression $Y = \alpha_0 + \alpha_1 X = u_i$ by OLS method and obtain the residuals ...

ii. If in fact the unrestricted regression $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 = u_i$ is the true regression, the ... obtained from estimated regression $Y = \alpha_0 + \alpha_1 X + u_i$ therefore ... related to the squared and cubed output terms, i.e., $X_i^2$ and $X_i^3$.

iii. This suggests that we regress the ... obtained in step ... on all the regressors including those in the restricted regression which in the present case means

$e = \alpha_0 + \alpha_1 X + ...$

where ... is an error term ... ... properties.

iv. For large sample size (large ... ) it is shown that $n$ (the sample size) times $R^2$ estimated from the regression equation of (3) follows a chi-square distribution with degrees of freedom ... equal to the number of restrictions imposed by the restricted regression ... the present example since the terms $X_i^3$ and $X_i^2$ are dropped from the model.

Symbolically we write $nR^2 \sim \chi_2^2$      (4) (no. of restrictions = 2)

## EXERCISE

21. The following table shows the values of expenditure on clothing $Y$, with expenditure $(X_1)$ and the price of clothing $(X_2)$

| | 1960 | 1961 | 1962 | 1963 | 1964 | 1965 | 1966 | 1967 | 1968 | 1969 |
|---|---|---|---|---|---|---|---|---|---|---|
| $X_2$ | 16 | 13 | 10 | 9 | 7 | 5 | 4 | 3 | 3.5 | 3 |
| $X_1$ | 15 | 20 | 30 | 42 | 50 | 54 | 63 | 72 | 85 | 96 |
| $Y$ | 1 4 | 4 3 | 5 | 6 | 7 | 9 | 8 | 10 | 12 | 4 |

(i) Estimate the model $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$

(ii) Estimate the model : $Y_i = \alpha_0 + \alpha_1 X_{1i} + v_i$

(iii) If $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$ is the true model, then examine the consequences on the regression parameters when $X_2$ is omitted from the model.

22. The following results were obtained from a sample of size 12

$\Sigma Y_i = 753$, $\Sigma Y_i^2 = 48{,}139$, $\Sigma X_{1i} Y_i = 40830$

$\Sigma X_{1i} = 643$, $\Sigma X_{1i}^2 = 34843$, $\Sigma X_{2i} Y_i = 6{,}796$

$\Sigma X_{2i} = 106$, $\Sigma X_{2i}^2 = 976$, $\Sigma X_{1i} X_{2i} = 5779$

(i) Estimate the model $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$

(ii) Estimate the model : $Y_i = \alpha_0 + \alpha_1 X_{1i} + v_i$

(iii) Examine the impact on the regression parameters when $X_2$ is omitted from the true model $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + u_i$

TABLE I (Contd.)

TABLE I (Contd.)

## TABLE II
### DISTRIBUTION OF STANDARD NORMAL VARIABLE
#### Values of $\tau_\alpha$

| $\alpha$ | 0.05 | 0.025 | 0.01 | 0.005 |
|---|---|---|---|---|
| $\tau_\alpha$ | 1.645 | 1.96 | 2.326 | 2.576 |

### TABLE III
### $\chi^2$ DISTRIBUTION
### VALUES OF $\chi^2$

For larger values of $n$, the quantity $\sqrt{2\chi^2} - \sqrt{2n-1}$ may be used as a standard normal variable.

*Abridged from Table 8 of *Biometrika Tables for Statisticians*, vol. I, with the kind permission of the Biometrika Trustees.

### TABLE IV
### $t$-DISTRIBUTION
### VALUES OF $t_{n,\alpha}$

Example

## TABLE VI
## THE DURBIN-WATSON d STATISTIC
### SIGNIFICANCE POINTS OF $d_L$ AND $d_U$ : 5%

| n | $k'=1$ $d_L$ | $d_U$ | $k'=2$ $d_L$ | $d_U$ | $k'=3$ $d_L$ | $d_U$ | $k'=4$ $d_L$ | $d_U$ | $k'=5$ $d_L$ | $d_U$ |
|---|---|---|---|---|---|---|---|---|---|---|
| 6 | 0.61 | 1.40 | — | — | — | — | — | — | — | — |
| 7 | 0.70 | 1.36 | 0.47 | 1.90 | — | — | — | — | — | — |
| 8 | 0.76 | 1.33 | 0.56 | 1.77 | 0.36 | 2.29 | — | — | — | — |
| 9 | 0.82 | 1.32 | 0.63 | 1.70 | 0.46 | 2.13 | 0.29 | 1.59 | — | — |
| 10 | 0.87 | 1.32 | 0.69 | 1.64 | 0.52 | 2.01 | 0.37 | 1.41 | 0.24 | 2.38 |
| 11 | 0.92 | 1.32 | 0.65 | 1.60 | 0.59 | 1.92 | 0.44 | 1.28 | 0.31 | 2.64 |
| 12 | 0.97 | 1.33 | 0.81 | 1.57 | 0.63 | 1.86 | 0.51 | 1.17 | 0.47 | 2.50 |
| 13 | 1.01 | 1.34 | 0.86 | 1.56 | 0.71 | 1.81 | 0.57 | 1.09 | 0.44 | 2.39 |
| 14 | 1.04 | 1.35 | 0.90 | 1.55 | 0.76 | 1.77 | 0.63 | 1.03 | 0.50 | 2.30 |
| 15 | 1.08 | 1.36 | 0.95 | 1.54 | 0.81 | 1.75 | 0.69 | 1.07 | 0.56 | 2.21 |
| 16 | 1.11 | 1.37 | 0.98 | 1.54 | 0.86 | 1.73 | 0.74 | 1.03 | 0.62 | 2.15 |
| 17 | 1.13 | 1.38 | 1.02 | 1.54 | 0.90 | 1.71 | 0.78 | 1.00 | 0.67 | 2.10 |
| 18 | 1.16 | 1.39 | 1.05 | 1.53 | 0.93 | 1.69 | 0.82 | 1.87 | 0.71 | 2.06 |
| 19 | 1.18 | 1.40 | 1.06 | 1.53 | 0.97 | 1.68 | 0.86 | 1.83 | 0.75 | 2.49 |
| 20 | 1.20 | 1.41 | 1.10 | 1.54 | 1.00 | 1.68 | 0.90 | 1.83 | 0.79 | 1.99 |
| 21 | 1.22 | 1.42 | 1.13 | 1.54 | 1.03 | 1.67 | 0.93 | 1.81 | 0.83 | 1.96 |
| 22 | 1.24 | 1.43 | 1.15 | 1.54 | 1.05 | 1.66 | 0.96 | 1.80 | 0.86 | 1.94 |
| 23 | 1.26 | 1.44 | 1.17 | 1.54 | 1.08 | 1.66 | 0.99 | 1.79 | 0.90 | 1.92 |
| 24 | 1.27 | 1.45 | 1.19 | 1.55 | 1.10 | 1.66 | 1.01 | 1.78 | 0.93 | 1.90 |
| 25 | 1.29 | 1.45 | 1.21 | 1.55 | 1.12 | 1.66 | 1.04 | 1.77 | 0.95 | 1.89 |
| 26 | 1.30 | 1.46 | 1.22 | 1.55 | 1.14 | 1.65 | 1.06 | 1.76 | 0.98 | 1.88 |
| 27 | 1.32 | 1.47 | 1.24 | 1.56 | 1.16 | 1.65 | 1.08 | 1.76 | 1.01 | 1.86 |
| 28 | 1.33 | 1.48 | 1.26 | 1.56 | 1.18 | 1.65 | 1.10 | 1.75 | 1.03 | 1.85 |
| 29 | 1.34 | 1.48 | 1.27 | 1.56 | 1.20 | 1.65 | 1.12 | 1.74 | 1.05 | 1.84 |
| 30 | 1.35 | 1.49 | 1.28 | 1.57 | 1.21 | 1.65 | 1.14 | 1.74 | 1.07 | 1.83 |
| 31 | 1.36 | 1.50 | 1.30 | 1.57 | 1.23 | 1.65 | 1.16 | 1.74 | 1.09 | 1.83 |
| 32 | 1.37 | 1.50 | 1.31 | 1.57 | 1.24 | 1.65 | 1.18 | 1.73 | 1.11 | 1.82 |
| 33 | 1.38 | 1.51 | 1.32 | 1.58 | 1.26 | 1.65 | 1.19 | 1.73 | 1.13 | 1.81 |
| 34 | 1.39 | 1.51 | 1.33 | 1.58 | 1.27 | 1.65 | 1.21 | 1.73 | 1.15 | 1.81 |
| 35 | 1.40 | 1.52 | 1.34 | 1.58 | 1.28 | 1.65 | 1.22 | 1.73 | 1.16 | 1.80 |
| 36 | 1.41 | 1.52 | 1.35 | 1.59 | 1.29 | 1.65 | 1.24 | 1.73 | 1.18 | 1.80 |
| 37 | 1.42 | 1.53 | 1.36 | 1.59 | 1.31 | 1.66 | 1.25 | 1.72 | 1.19 | 1.80 |
| 38 | 1.43 | 1.54 | 1.37 | 1.59 | 1.32 | 1.66 | 1.26 | 1.72 | 1.21 | 1.79 |
| 39 | 1.43 | 1.54 | 1.38 | 1.60 | 1.33 | 1.66 | 1.27 | 1.72 | 1.22 | 1.79 |
| 40 | 1.44 | 1.54 | 1.39 | 1.60 | 1.34 | 1.66 | 1.29 | 1.72 | 1.23 | 1.79 |
| 45 | 1.47 | 1.57 | 1.43 | 1.62 | 1.38 | 1.67 | 1.34 | 1.72 | 1.29 | 1.78 |
| 50 | 1.50 | 1.59 | 1.46 | 1.63 | 1.42 | 1.67 | 1.38 | 1.72 | 1.34 | 1.77 |
| 55 | 1.53 | 1.60 | 1.49 | 1.64 | 1.45 | 1.68 | 1.41 | 1.72 | 1.38 | 1.77 |
| 60 | 1.55 | 1.62 | 1.51 | 1.65 | 1.48 | 1.69 | 1.44 | 1.73 | 1.41 | 1.77 |
| 65 | 1.57 | 1.63 | 1.54 | 1.66 | 1.50 | 1.70 | 1.47 | 1.73 | 1.44 | 1.77 |
| 70 | 1.58 | 1.64 | 1.55 | 1.67 | 1.52 | 1.70 | 1.49 | 1.74 | 1.46 | 1.77 |
| 75 | 1.60 | 1.65 | 1.57 | 1.68 | 1.54 | 1.71 | 1.51 | 1.74 | 1.49 | 1.77 |
| 80 | 1.61 | 1.66 | 1.59 | 1.69 | 1.56 | 1.72 | 1.53 | 1.74 | 1.51 | 1.77 |
| 85 | 1.62 | 1.67 | 1.60 | 1.70 | 1.57 | 1.72 | 1.55 | 1.75 | 1.52 | 1.77 |
| 90 | 1.63 | 1.68 | 1.61 | 1.70 | 1.59 | 1.73 | 1.57 | 1.75 | 1.54 | 1.78 |
| 95 | 1.64 | 1.69 | 1.62 | 1.71 | 1.60 | 1.73 | 1.58 | 1.75 | 1.56 | 1.78 |
| 100 | 1.65 | 1.69 | 1.63 | 1.72 | 1.61 | 1.74 | 1.59 | 1.76 | 1.57 | 1.78 |
| 150 | 1.72 | 1.74 | 1.70 | 1.76 | 1.69 | 1.77 | 1.87 | 1.78 | 1.66 | 1.80 |
| 200 | 1.75 | 1.77 | 1.74 | 1.78 | 1.73 | 1.79 | 1.72 | 1.81 | 1.71 | 1.82 |

Note : $k'$ = Number of explanatory variables excluding the constant.